# Synchronization Cache Coherence

Chris Rossbach

CS378H

# Today

- Questions?
- Administrivia
  - Lab 1 due soon
- Material for the day
  - Cache coherence
  - Lock implementation
  - Blocking synchronization

- Acknowledgements
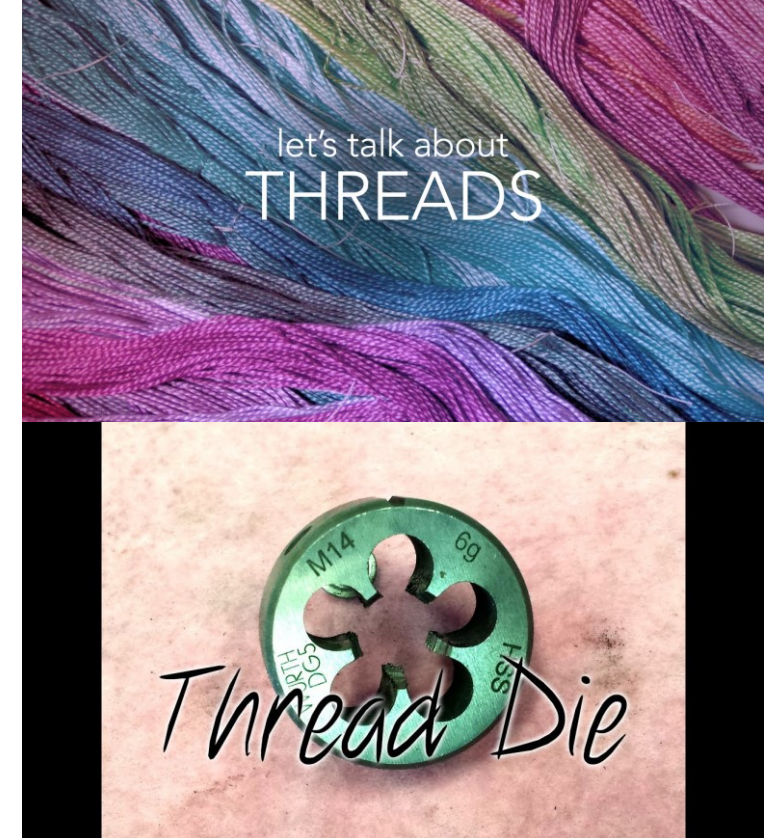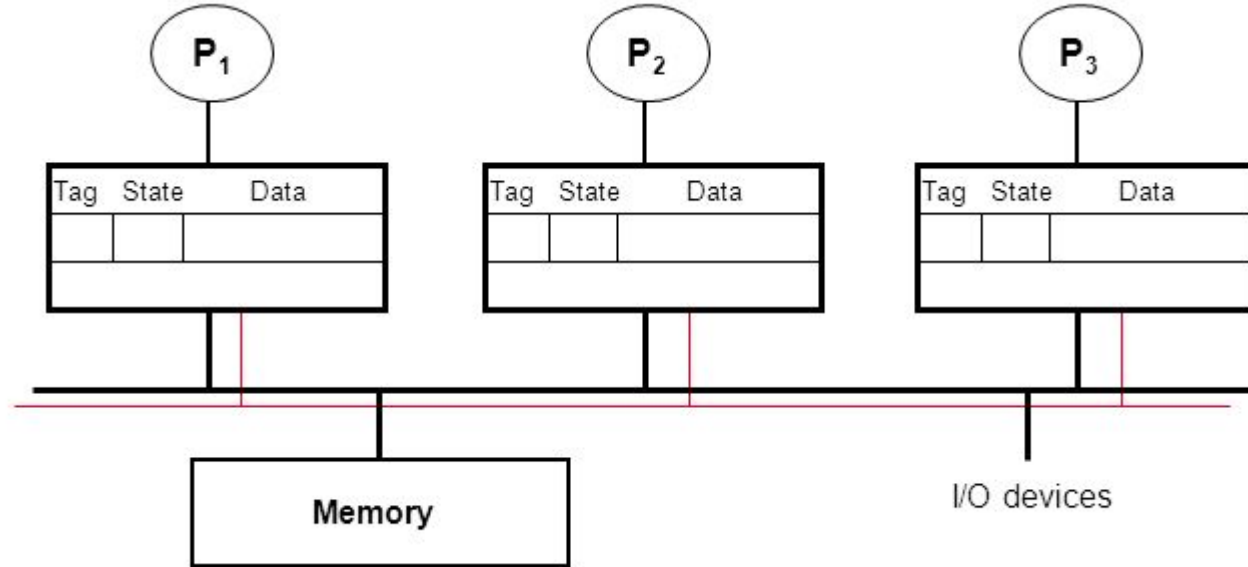  - Thanks to Gadi Taubenfield: I borrowed from some of his slides on barriers

# Today

- Questions?
- Administrivia
  - Lab 1 due soon
- Material for the day
  - Cache coherence
  - Lock implementation
  - Blocking synchronization

- Acknowledgements
  - Thanks to Gadi Taubenfield: I borrowed from some of his slides on barriers

let's talk about
THREADS

# Today

- Questions?
- Administrivia
  - Lab 1 due soon
- Material for the day
  - Cache coherence
  - Lock implementation
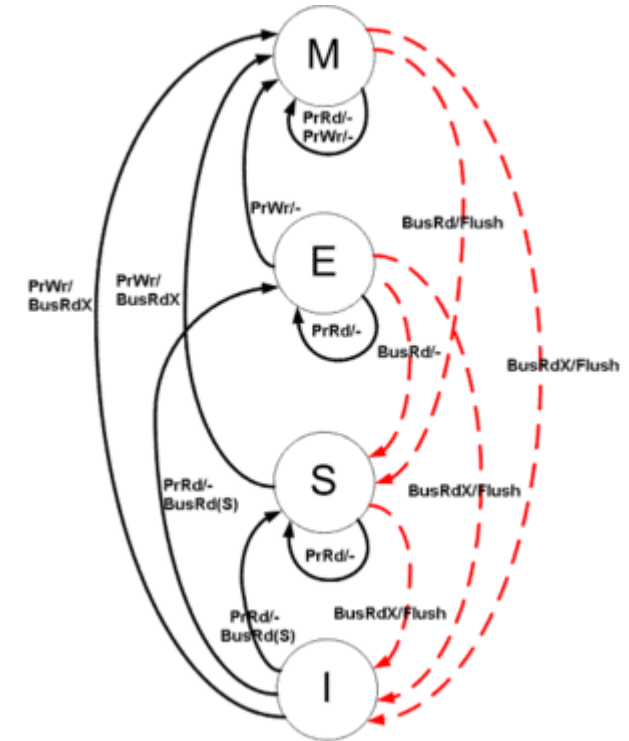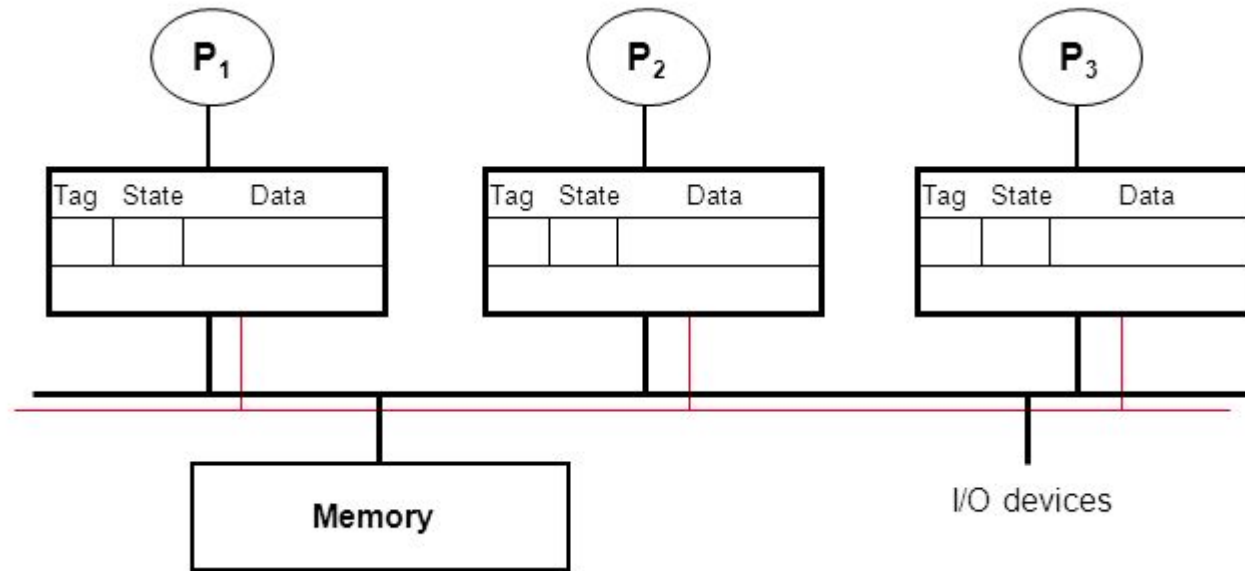  - Blocking synchronization

# Faux Quiz (answer any 2, 5 min)

- What is the difference between spinning/busy-wait and blocking synchronization?
- Can you write shared memory parallel applications using single-threaded processes only?
- How do you choose between spinlock/mutex on a multi-processor?
- Define the states of the MESI protocol. Is the E state necessary? Why or why not?
- What is bus locking?
- What is the difference between Mesa and Hoare monitors?
- Why recheck the condition on wakeup from a monitor wait?
- How can you build barriers with spinlocks?
- How can you build barriers with monitors?
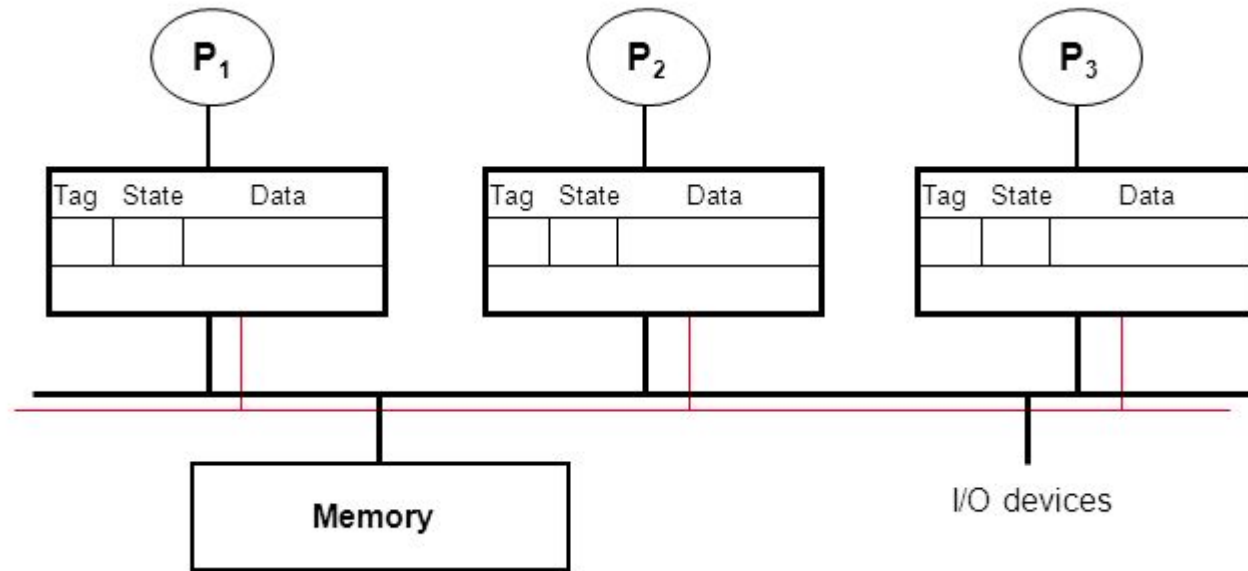- What is the difference between a mutex and a semaphore?

# Review: Basic MESI Cache Coherence
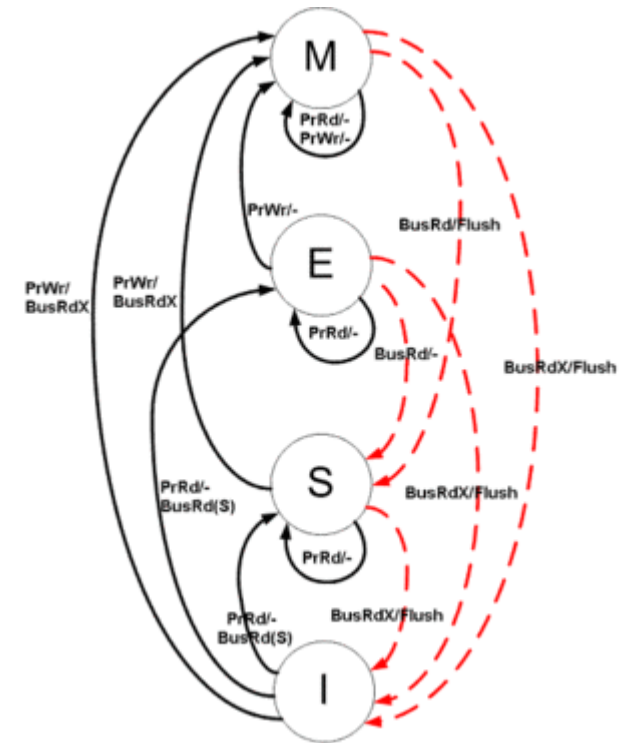
# Review: Basic MESI Cache Coherence

# Review: Basic MESI Cache Coherence



Each cache line has a state (M, E, S, I)

# Review: Basic MESI Cache Coherence



Each cache line has a state (M, E, S, I)
- Processors "snoop" bus to maintain states

# Review: Basic MESI Cache Coherence



Each cache line has a state (M, E, S, I)
- Processors "snoop" bus to maintain states
- Initially → 'I' → Invalid

**INVALID**

# Review: Basic MESI Cache Coherence



Each cache line has a state (M, E, S, I)
- Processors "snoop" bus to maintain states
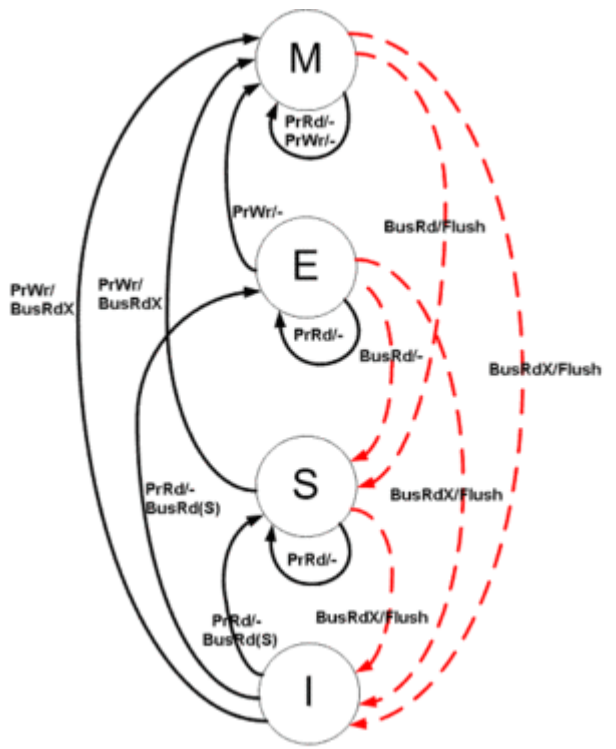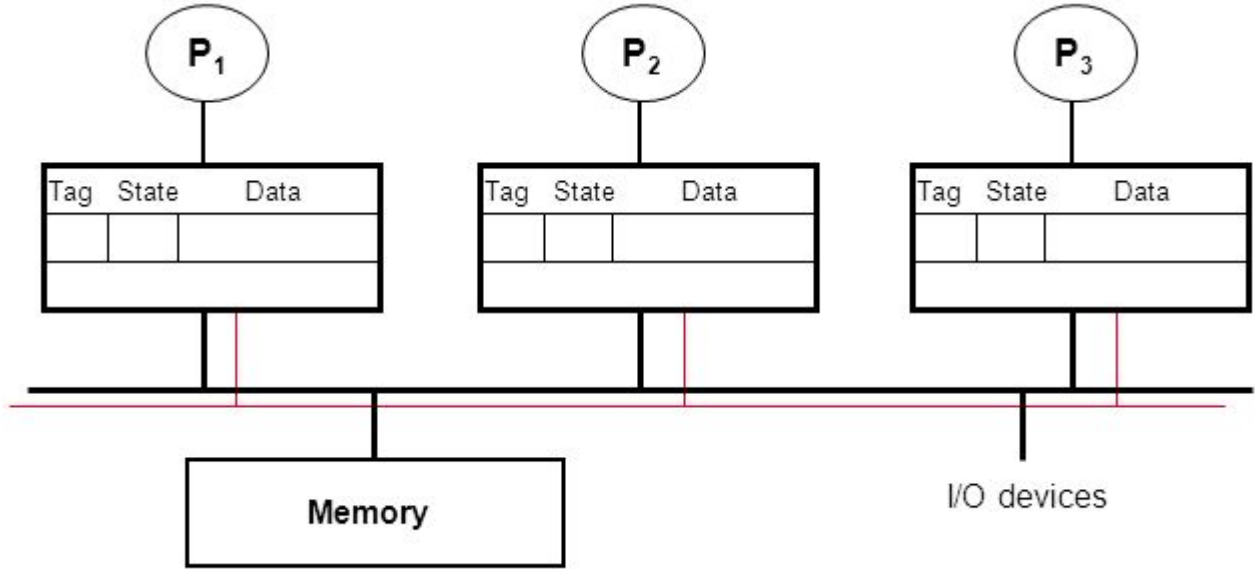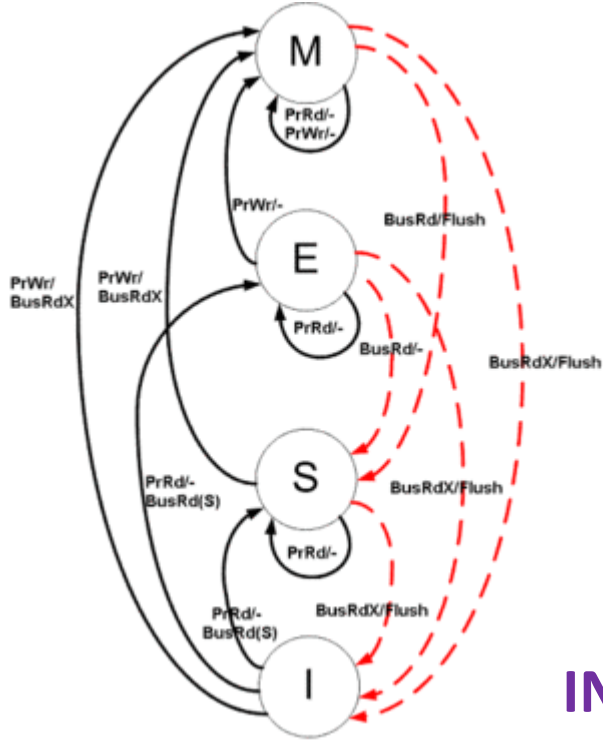- Initially → 'I' → Invalid
- Read one → 'E' → exclusive

# Review: Basic MESI Cache Coherence



Each cache line has a state (M, E, S, I)
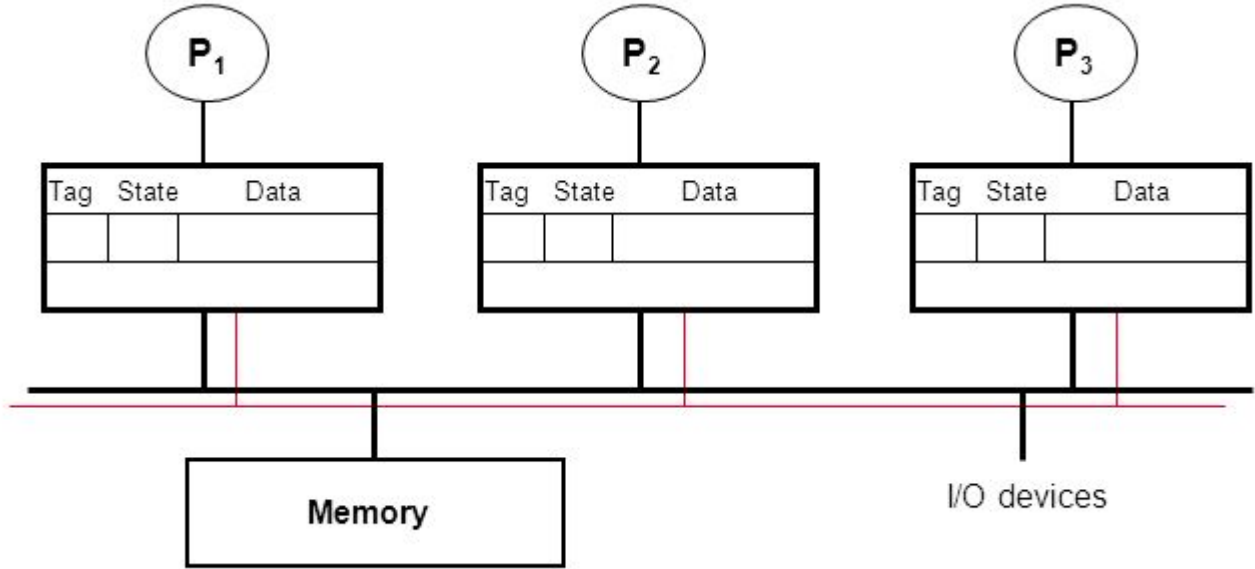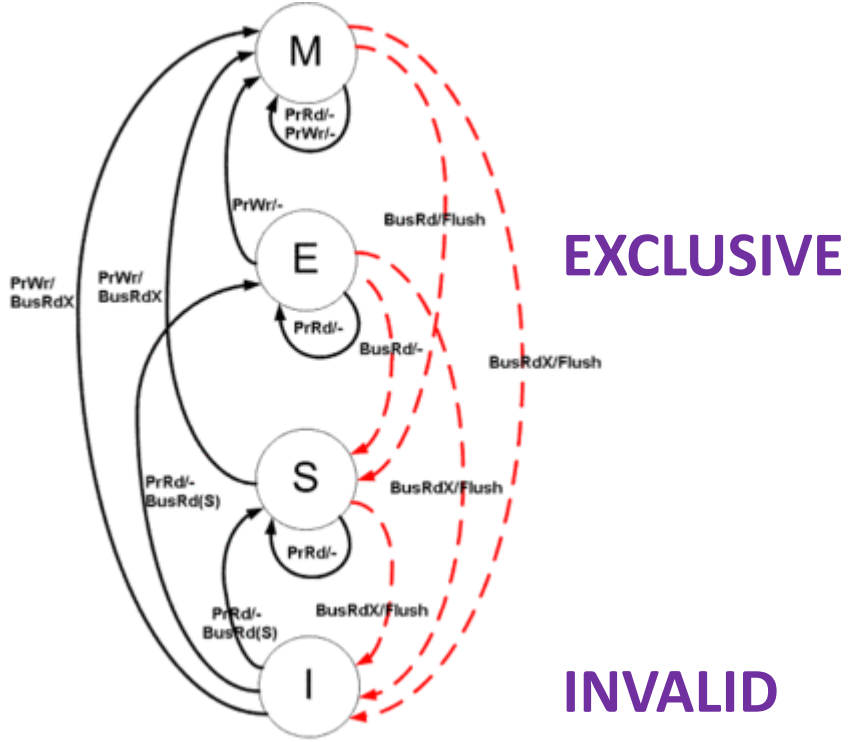- Processors "snoop" bus to maintain states
- Initially → 'I' → Invalid
- Read one → 'E' → exclusive
- Reads → 'S' → multiple copies possible

EXCLUSIVE

SHARED

INVALID

# Review: Basic MESI Cache Coherence



**EXCLUSIVE**

**SHARED**

**INVALID**

Each cache line has a state (M, E, S, I)

- Processors "snoop" bus to maintain states
- Initially → 'I' → Invalid
- Read one → 'E' → exclusive
- Reads → 'S' → multiple copies possible
- Write → 'M' → single copy → lots of cache coherence traffic

# Review: Basic MESI Cache Coherence



MODIFIED

EXCLUSIVE
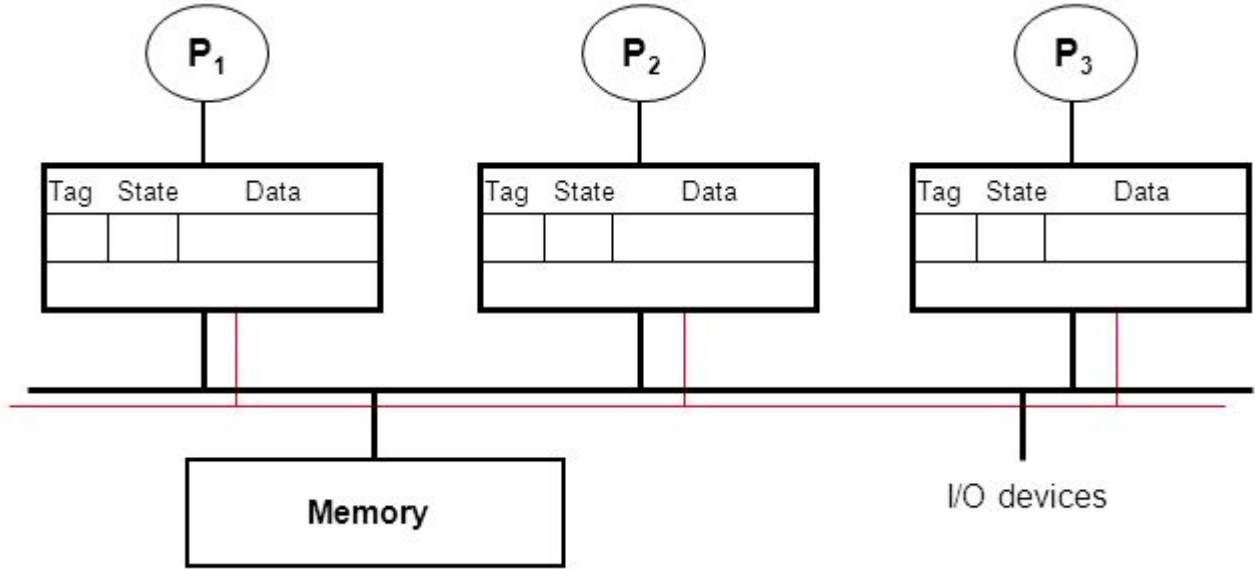
SHARED

INVALID
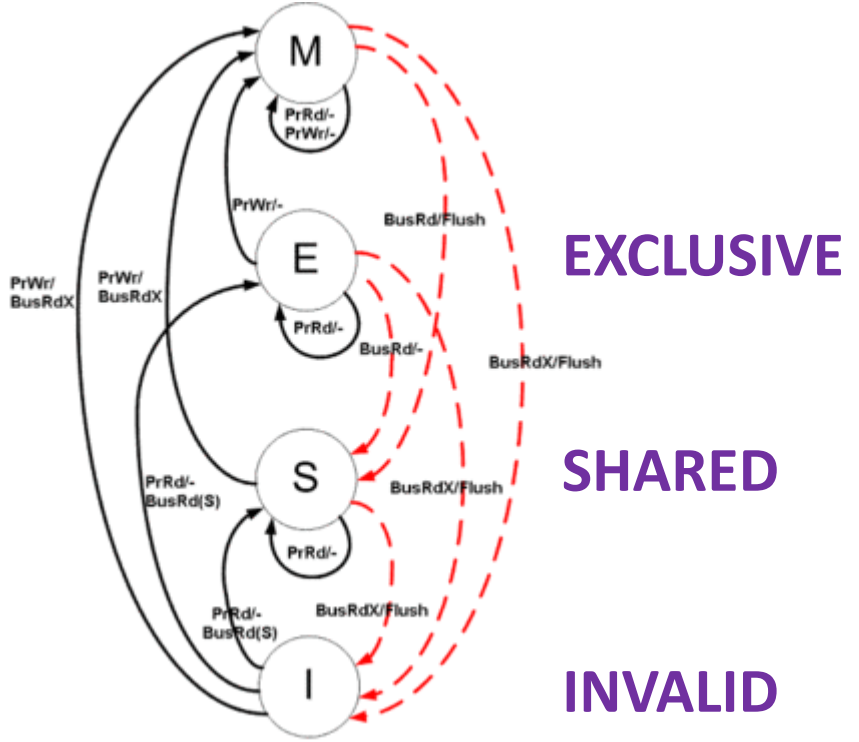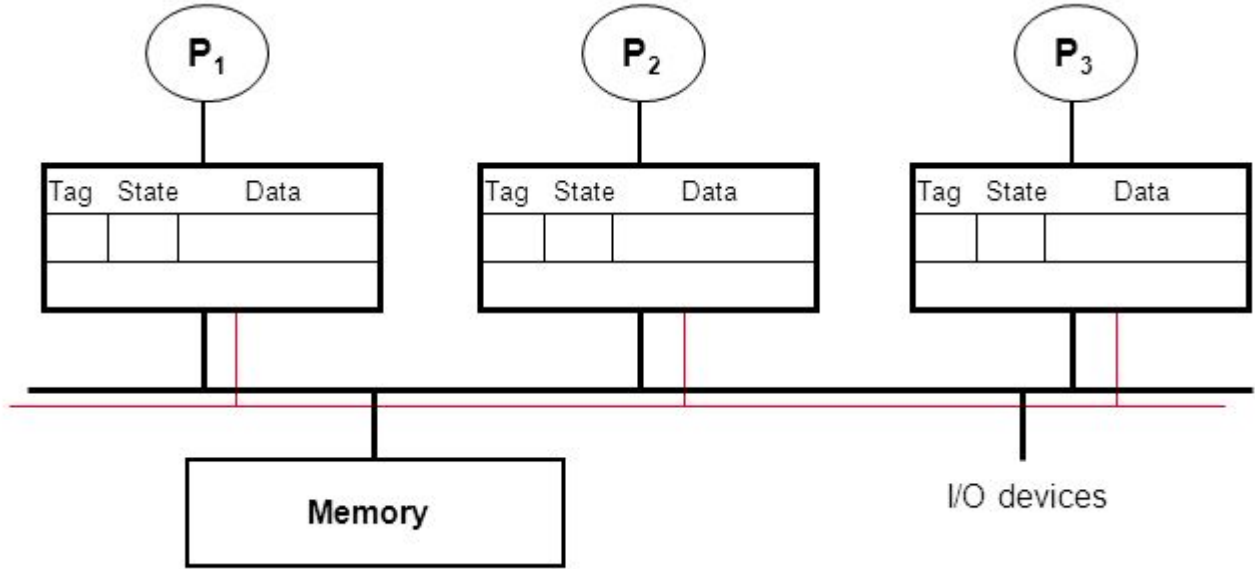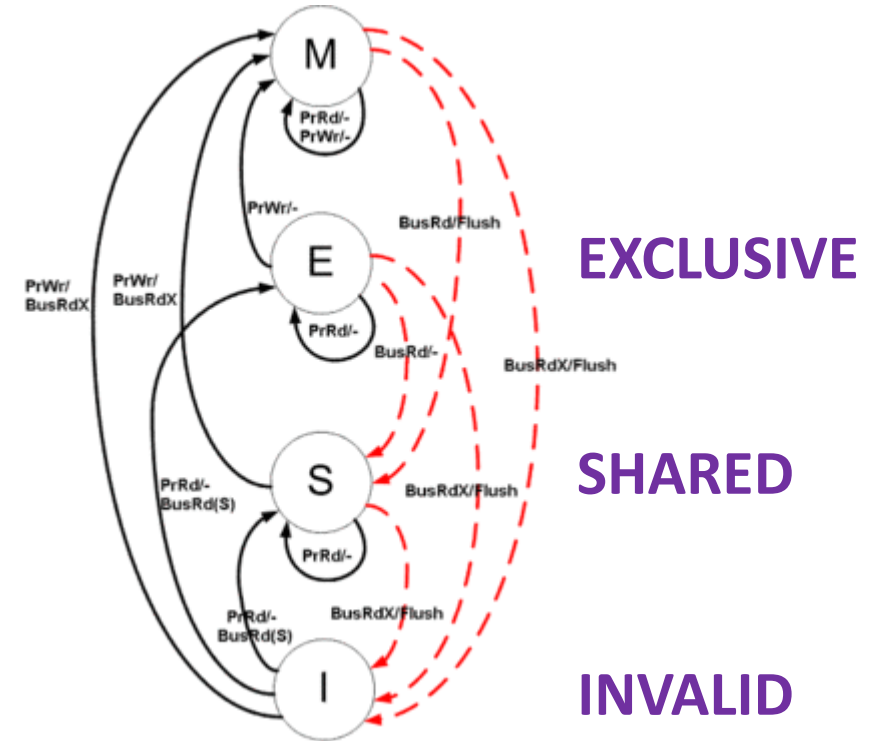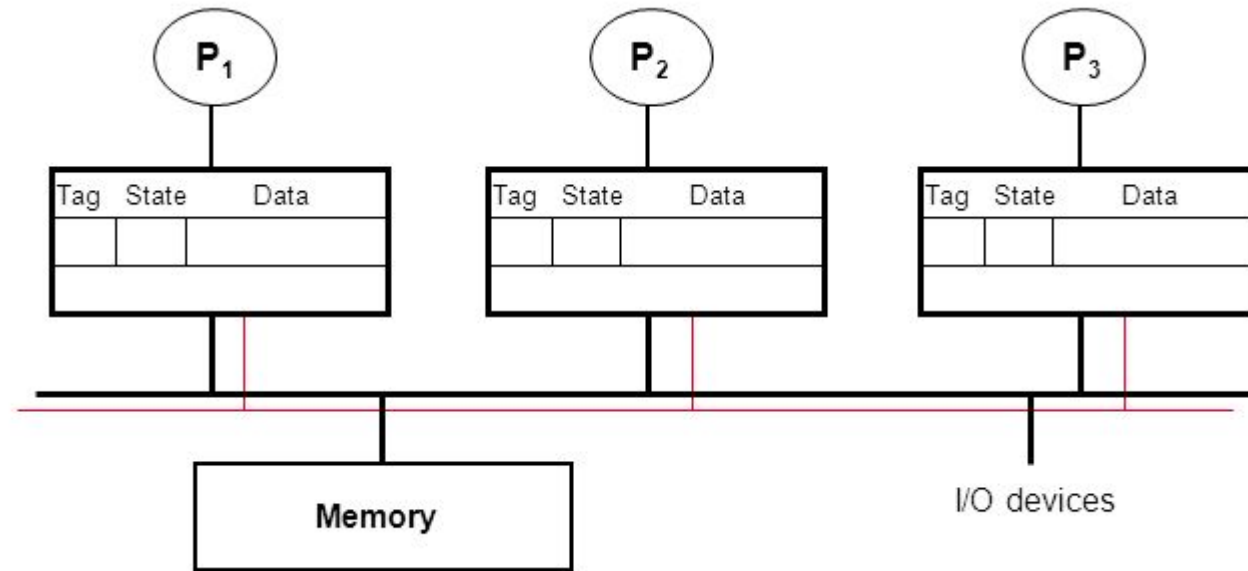
Each cache line has a state (M, E, S, I)

- Processors "snoop" bus to maintain states
- Initially → 'I' → Invalid
- Read one → 'E' → exclusive
- Reads → 'S' → multiple copies possible
- Write → 'M' → single copy → lots of cache coherence traffic

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED

INVALID

cache      cache      cache

lock:     lock:

Memory
lock: 0

I/O devices

P1

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED

INVALID

P1

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

P2

```
// (straw-person lock impl)
// Initilly, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

# Cache Coherence Action Zone



| | |
|---|---|
| M | **MODIFIED** |
| E | **EXCLUSIVE** |
| S | **SHARED** |
| I | **INVALID** |

P1

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED

INVALID

cache — P1

```
Tag  State    Data
      E  lock: 0
```

cache — P2

```
Tag  State    Data
      I  lock:
```

cache — P3

```
Tag  State    Data
```

Memory
lock: 0

I/O devices

**P1**

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

**P2**

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED
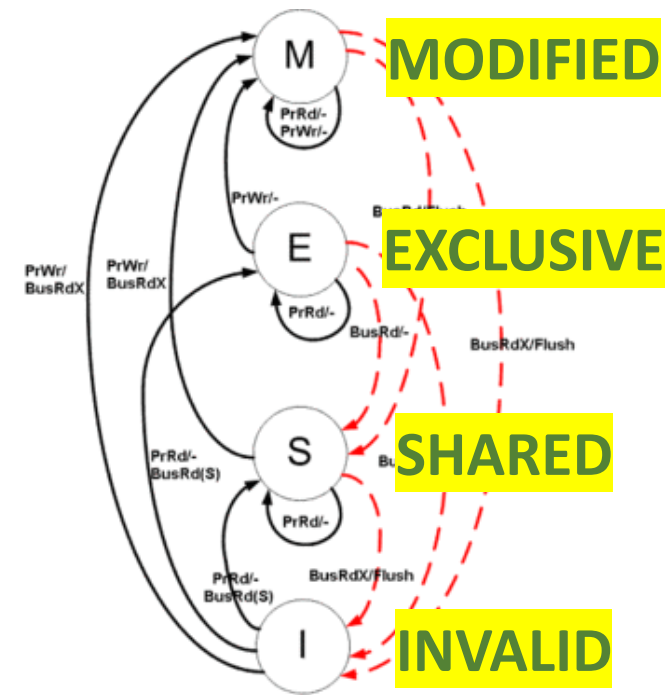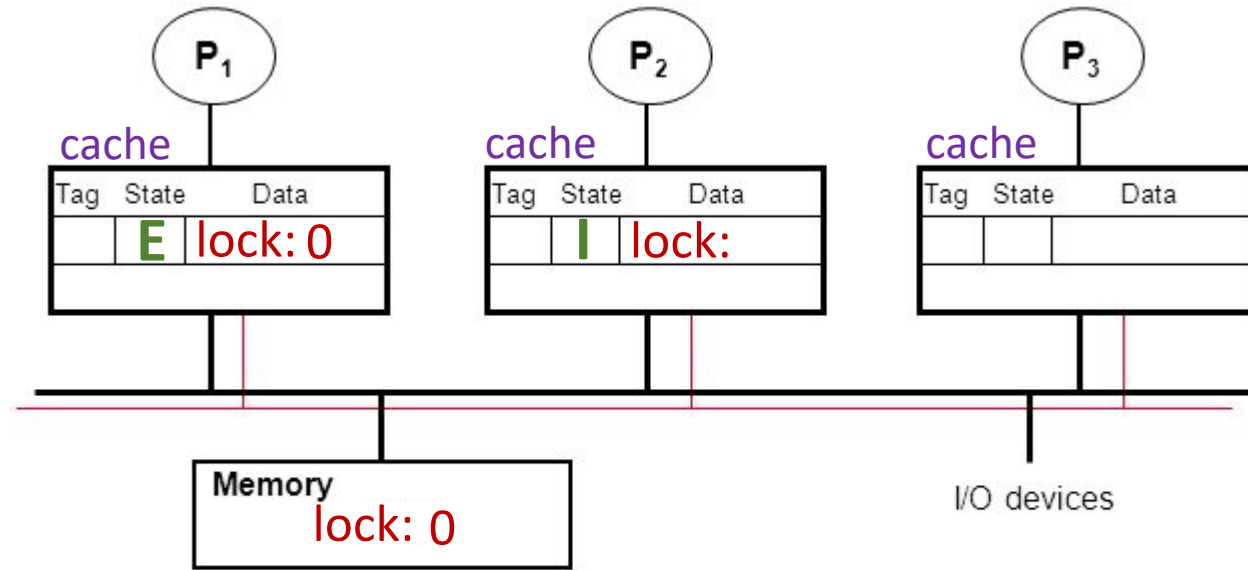
INVALID

P1            P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED

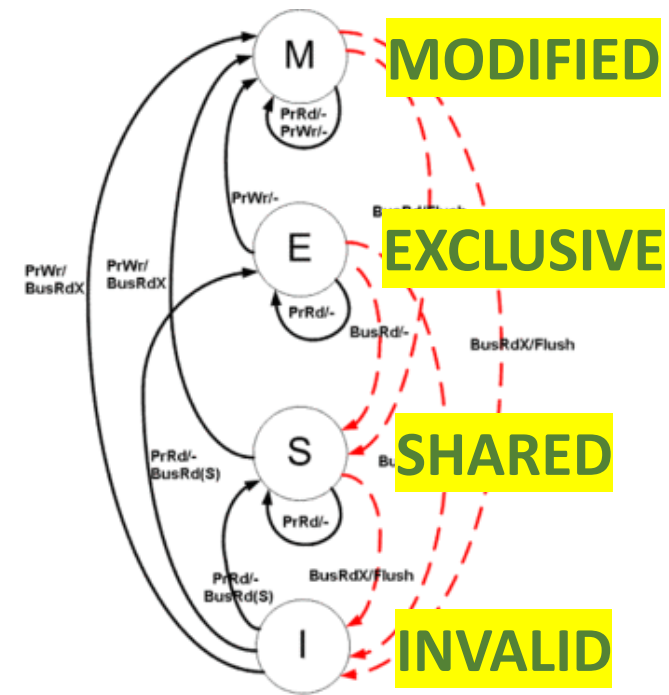EXCLUSIVE

SHARED

INVALID

P1

P2
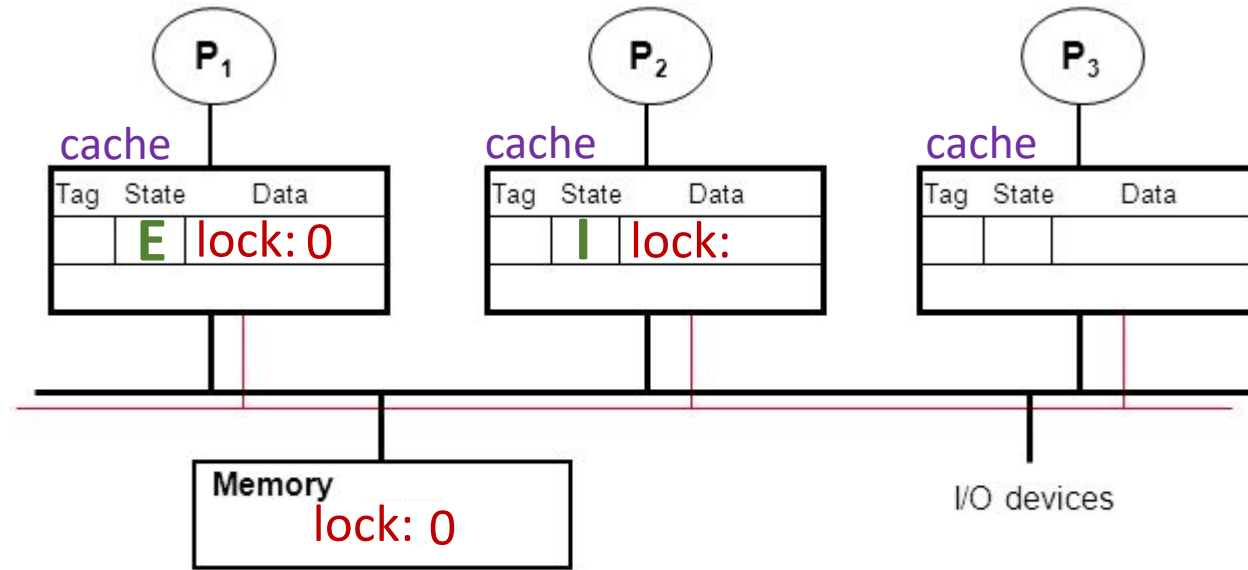
```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

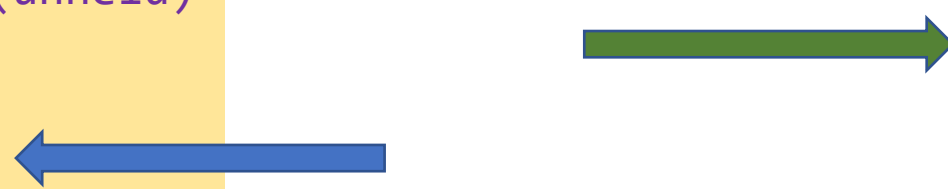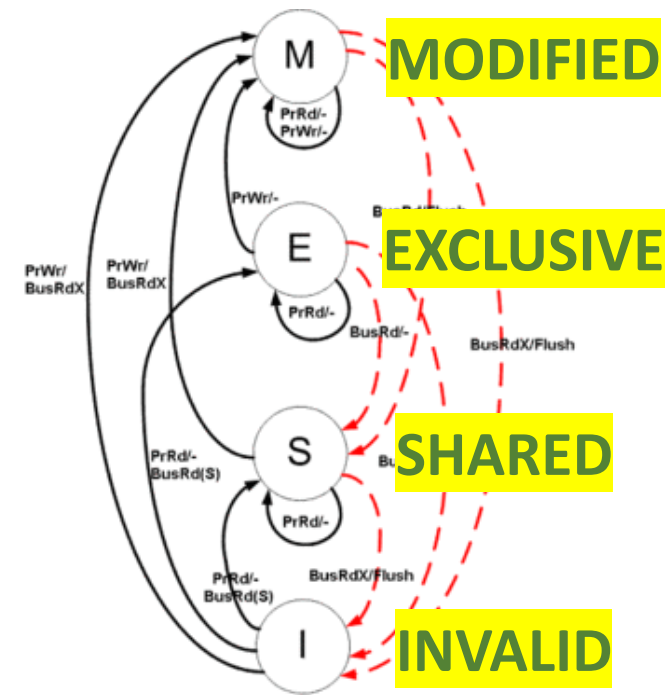# Cache Coherence Action Zone



P1

P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1

}
```

```
// (straw-person lock impl)
// Inilily, lock == 0 (unheld)
lock() {
try:   load lock, R0
          test R0
          bnz try
          store lock, 1

}
```

# Cache Coherence Action Zone



MODIFIED
EXCLUSIVE
SHARED
INVALID

P1

P2

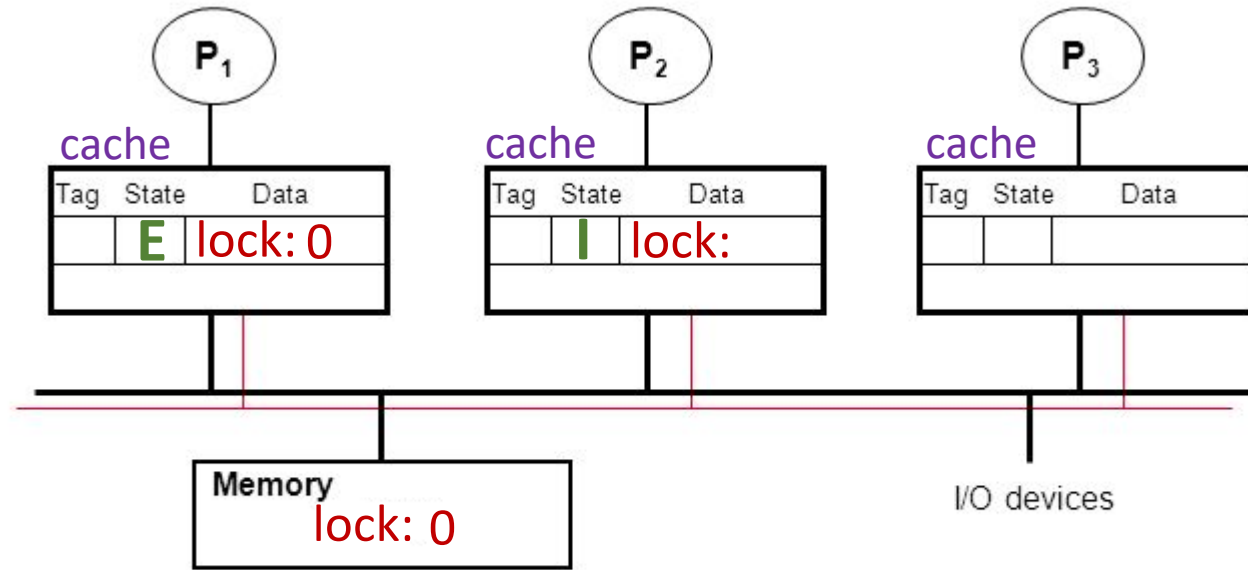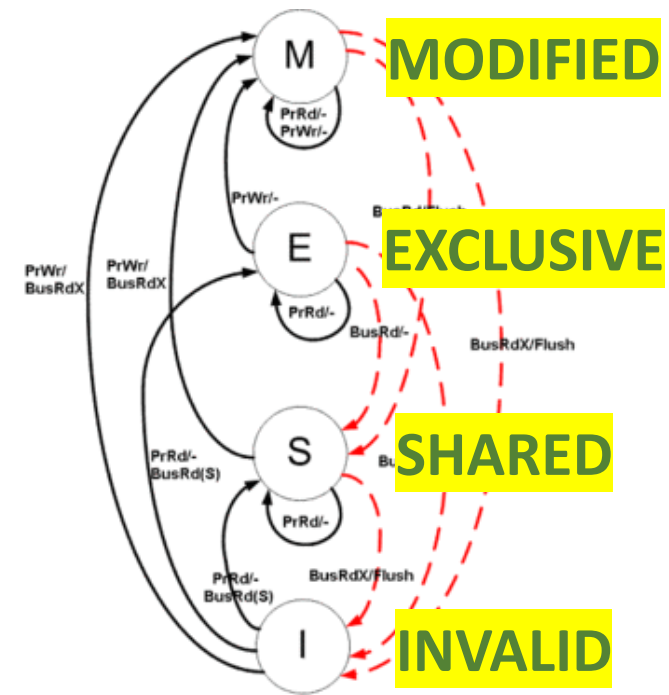```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED
EXCLUSIVE
SHARED
INVALID

cache

| Tag | State | Data |
|-----|-------|------|
|     | S     | lock:1 |
|     |       |      |

cache

| Tag | State | Data |
|-----|-------|------|
|     | S     | lock: 1 |
|     |       |      |

cache

| Tag | State | Data |
|-----|-------|------|
|     |       |      |
|     |       |      |

Memory
lock: 1

I/O devices

P1

P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```
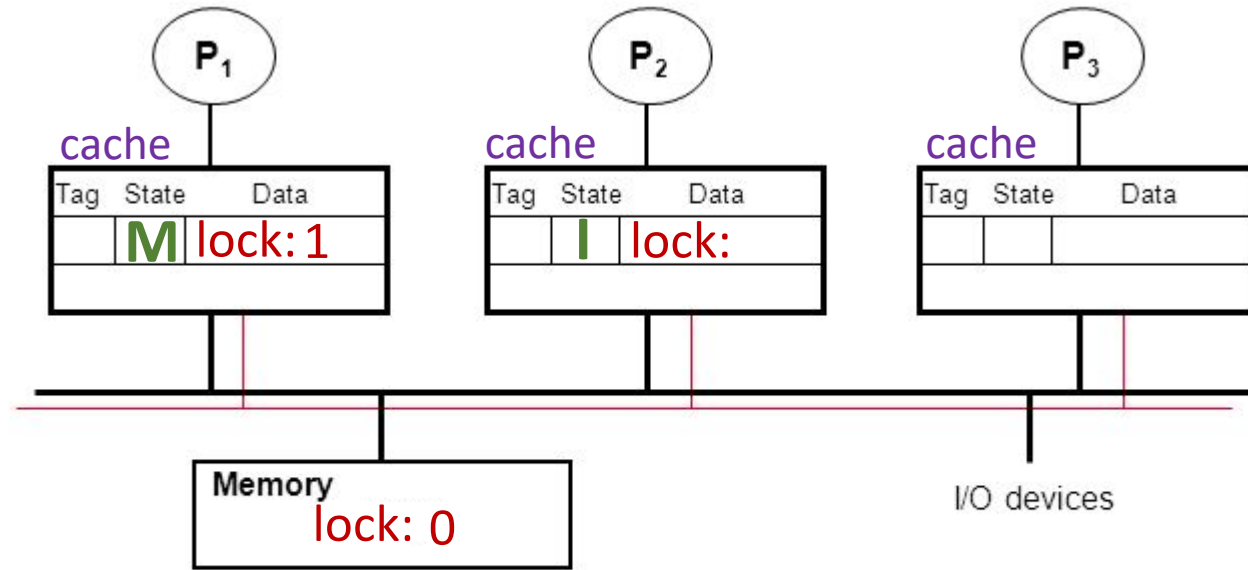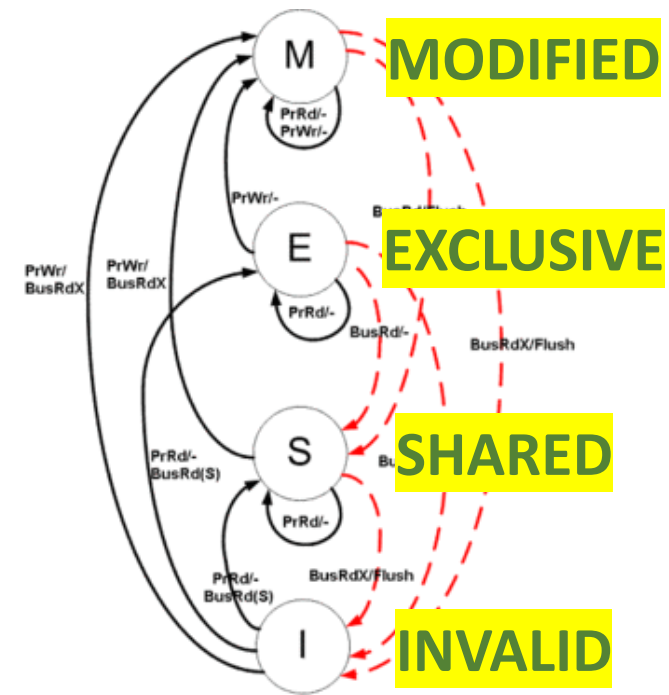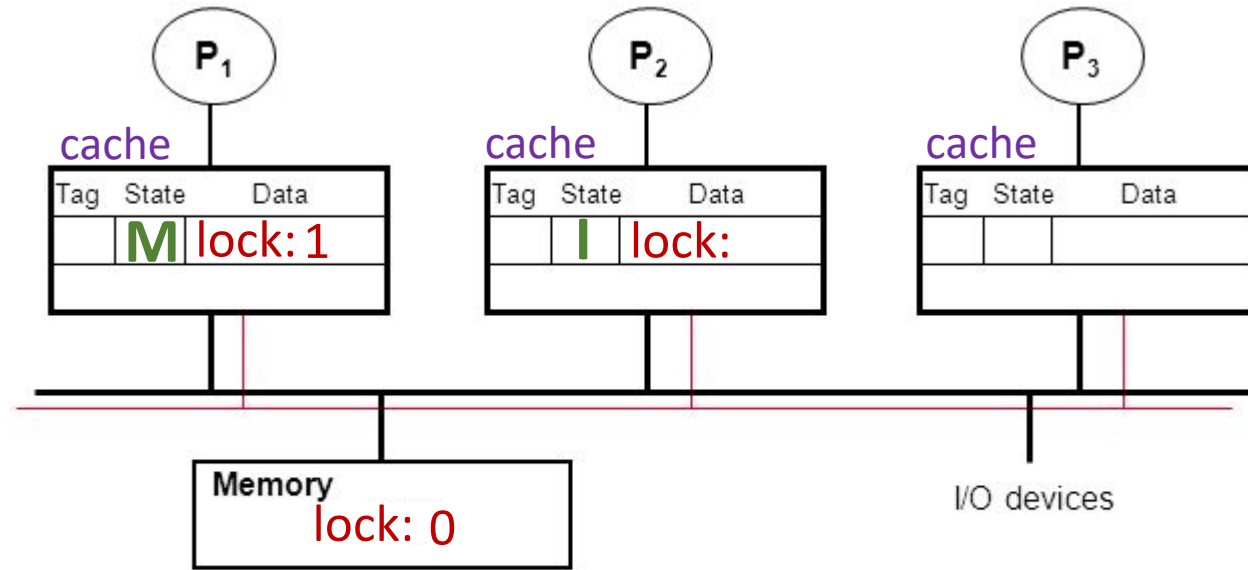
```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1
}
```

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED

INVALID

cache | cache | cache

Tag State Data — S lock:1

Tag State Data — S lock: 1

Tag State Data

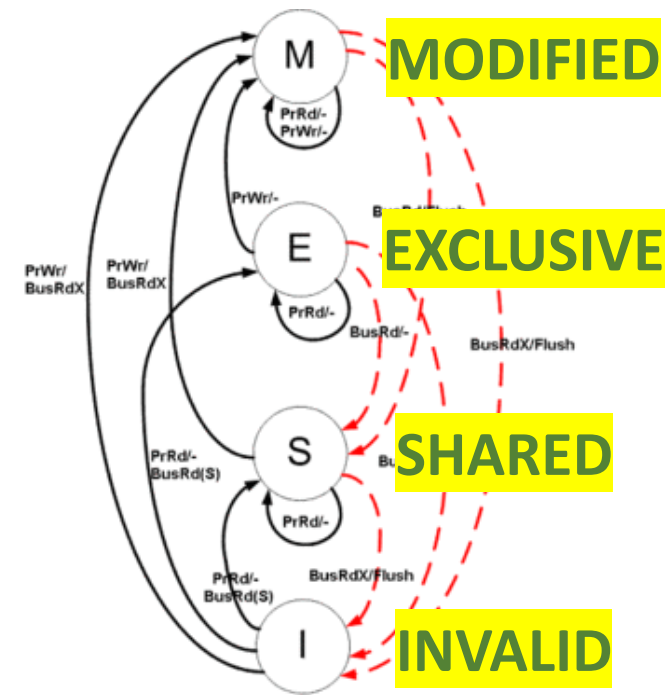Memory — lock: 1

I/O devices

P1

P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1

}
```

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
        test R0
        bnz try
        store lock, 1

}
```

# Cache Coherence Action Zone



MODIFIED

EXCLUSIVE

SHARED

INVALID

P1

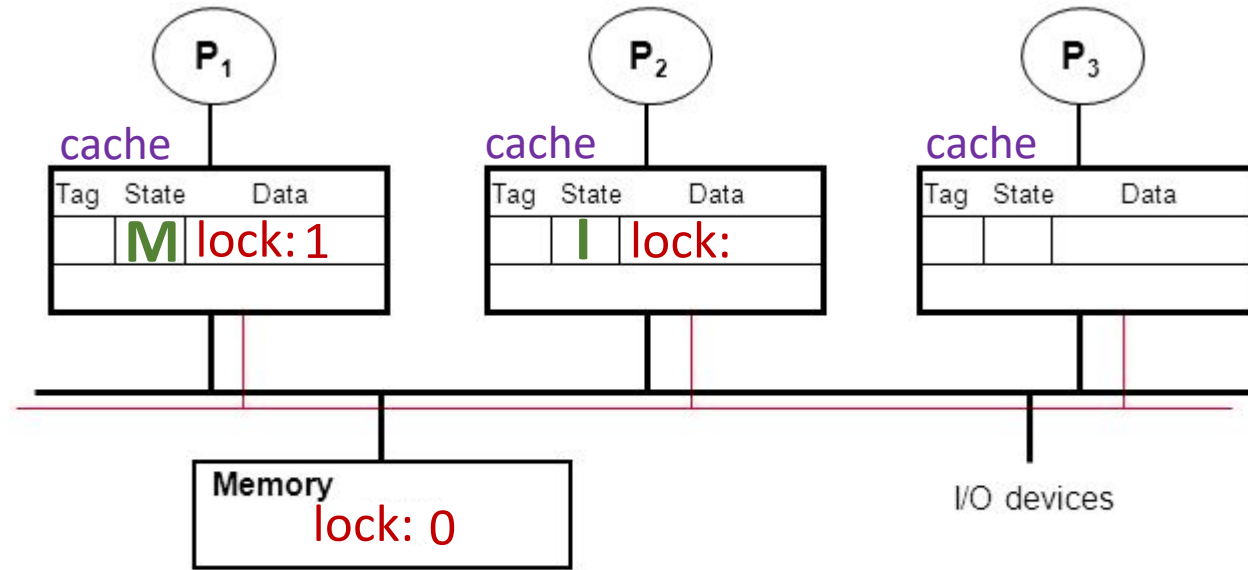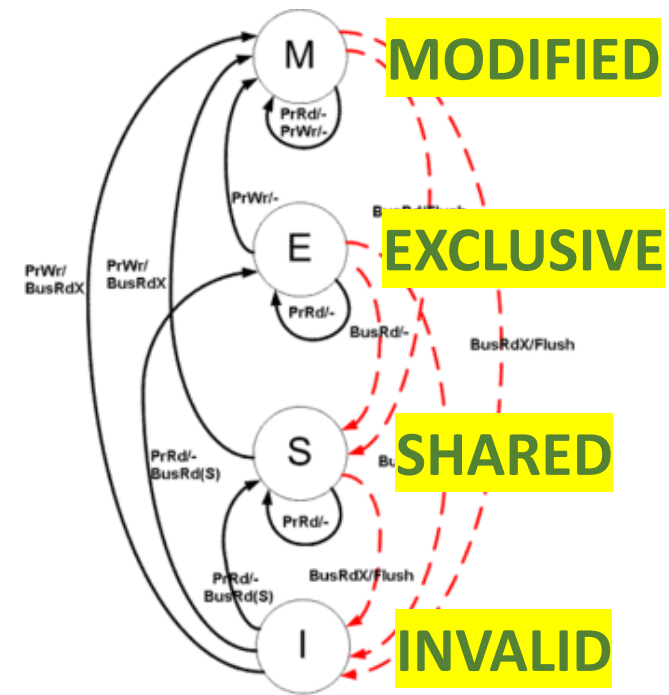P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:  load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

```
// (straw-person lock impl)
// Initilly, lock == 0 (unheld)
lock() {
try:  load lock, R0
          test R0
          bnz try
          store lock, 1
}
```

SAFE!

# Cache Coherence Action

**WAIT! Is E necessary?**



MODIFIED
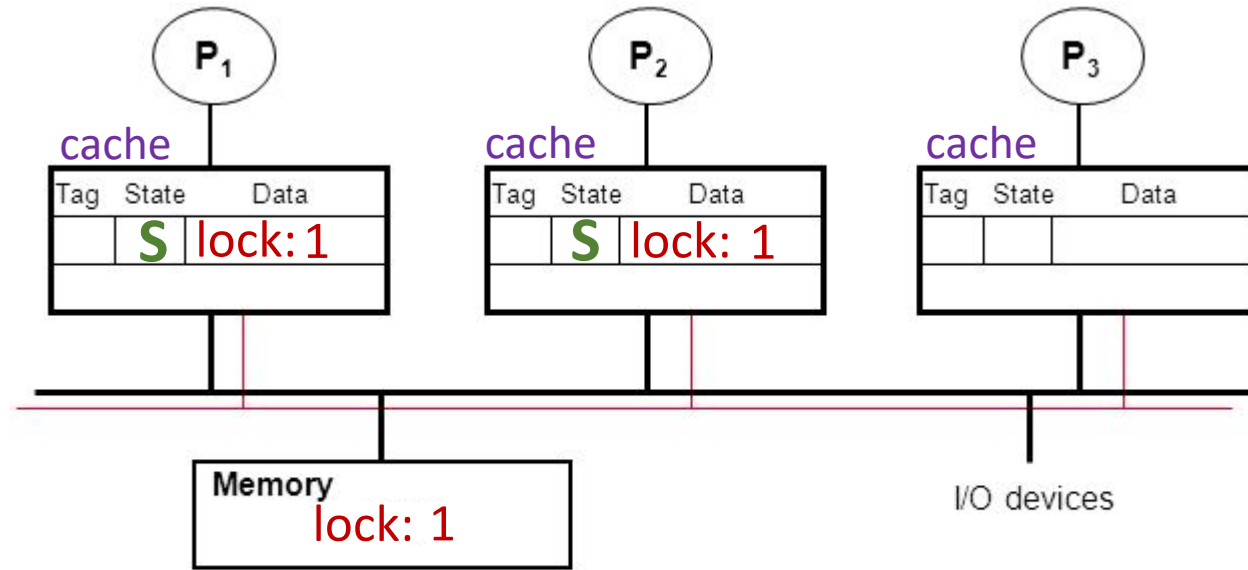EXCLUSIVE
SHARED
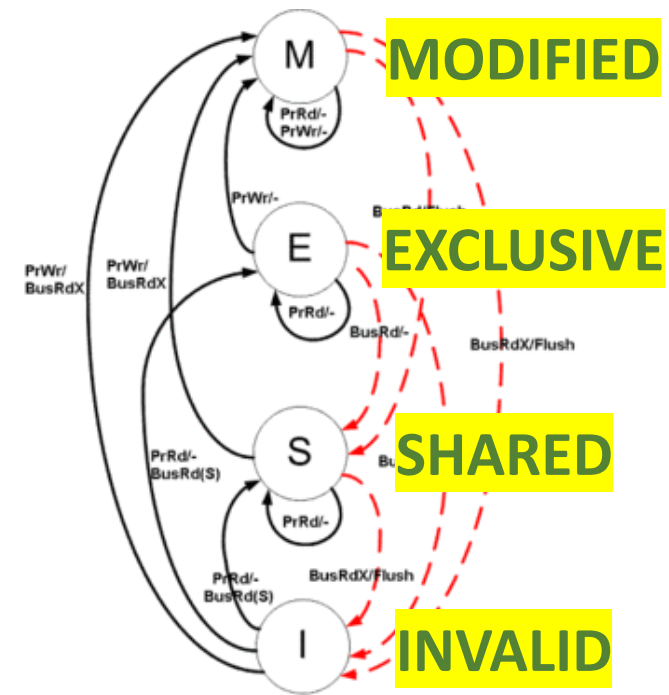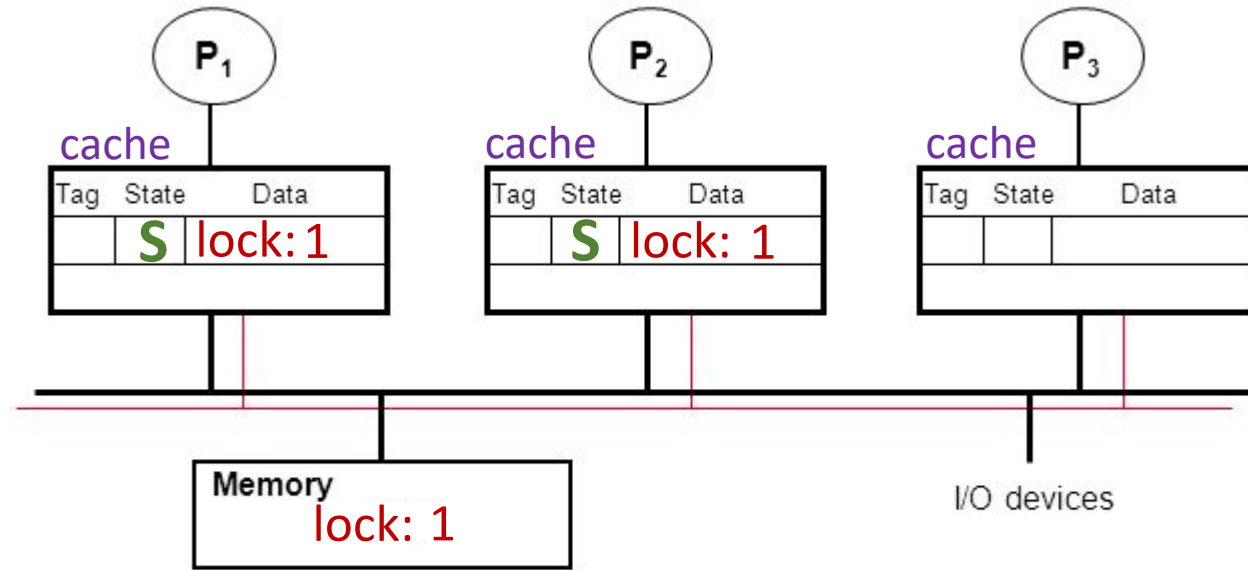INVALID

P1

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
       test R0
       bnz try
       store lock, 1
}
```
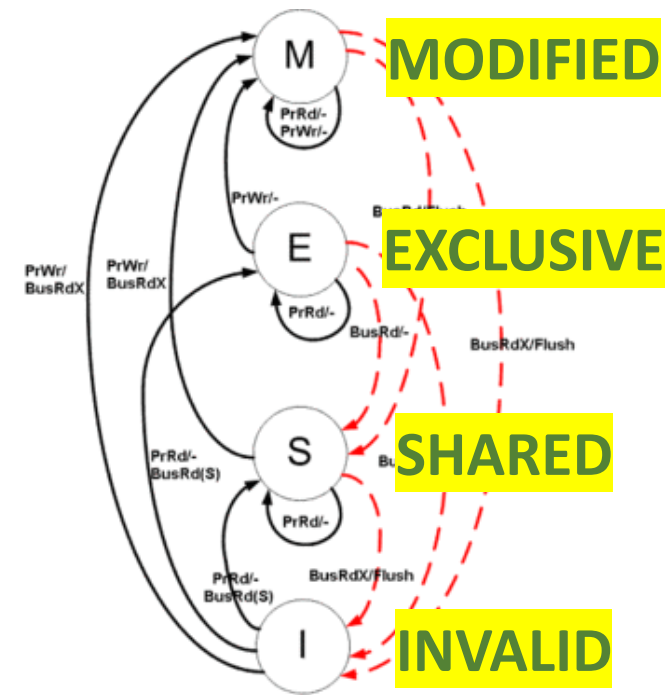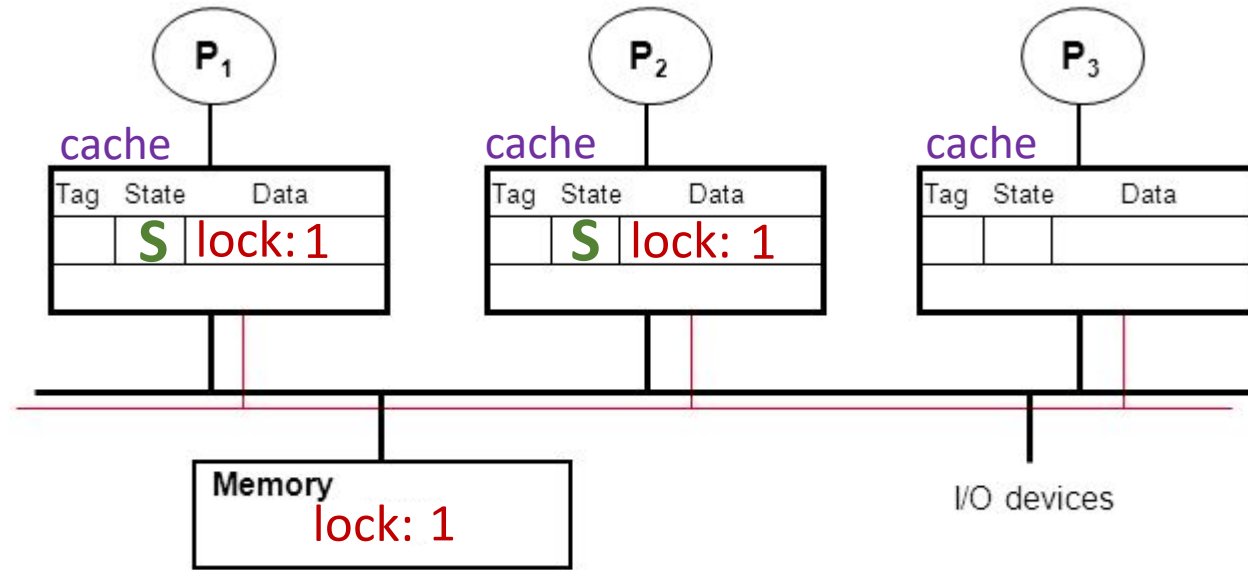
P2

```
// (straw-person lock impl)
// Initially, lock == 0 (unheld)
lock() {
try:   load lock, R0
       test R0
       bnz try
       store lock, 1
}
```

SAFE!

# Other Coherence Protocols: MSI

# Other Coherence Protocols: MOESI

# Other Coherence Protocols: FRMSI

# HW Support for RMW: LL-SC

| LLSC: load-linked store-conditional |
|---|
| PPC, Alpha, MIPS |

```
LL(addr, val) {          bool SC(addr, val) {
 link(addr);                 if(link-ok(addr)) {
 return *addr;                  *addr = val;
}                               return true;
                             }
                             return false;
                          }
```

- load-linked is a load that is "linked" to a subsequent store-conditional
- Store-conditional only succeeds if value from linked-load is unchanged

# HW Support for RMW: LL-SC

**LLSC: load-linked store-conditional**

PPC, Alpha, MIPS

```
LL(addr, val) {          bool SC(addr, val) {
 link(addr);                if(link-ok(addr)) {
 return *addr;                 *addr = val;
}                              return true;
                            }
                            return false;
                          }
```

```
void LLSC_lock(lock) {
   while(1) {
      old = load-linked(lock);
      if(old == 0 && store-cond(lock, 1))
         return;
   }
}
```

- load-linked is a load that is "linked" to a subsequent store-conditional
- Store-conditional only succeeds if value from linked-load is unchanged

# HW Support for RMW: LL-SC

**LLSC: load-linked store-conditional**

PPC, Alpha, MIPS

```
LL(addr, val) {          bool SC(addr, val) {
 link(addr);                if(link-ok(addr)) {
 return *addr;                  *addr = val;
}                              return true;
                             }
                             return false;
                           }
```

```
void LLSC_lock(lock) {
  while(1) {
    old = load-linked(lock);
    if(old == 0 && store-cond(lock, 1))
      return;
  }
}
```

- load-linked is a load that is "linked" to a subsequent store-conditional
- Store-conditional only succeeds if value from linked-load is unchanged

# LLSC Lock Action Zone II



```
_____P1_____
lock(lock) {
    while(1) {
        old = ll(lock);
        if(old == 0)
            if(sc(lock, 1))
                return;
    }
}
```

```
_____P2_____
lock(lock) {
    while(1) {
        old = ll(lock);
        if(old == 0)
            if(sc(lock, 1))
                return;
    }
}
```

# LLSC Lock Action Zone II



P1:
lock:

P2:
lock: S[L] 0

lock: 0

```
_____P1_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

```
_____P2_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

# LLSC Lock Action Zone II



State Data

lock:

State Data

lock: S[L] 0

lock: 0

```
_____P1_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

```
_____P2_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

# LLSC Lock Action Zone II



```
_____P1_____
lock(lock) {
   while(1) {
      old = ll(lock);
      if(old == 0)
         if(sc(lock, 1))
            return;
   }
}
```

```
_____P2_____
lock(lock) {
   while(1) {
      old = ll(lock);
      if(old == 0)
         if(sc(lock, 1))
            return;
   }
}
```

# LLSC Lock Action Zone II



P1 State Data
lock: S[L] 0

P2 State Data
lock: S[L] 0

lock: 0

```
_____P1_____
lock(lock) {
    while(1) {
        old = ll(lock);
        if(old == 0)
            if(sc(lock, 1))
                return;
    }
}
```

```
_____P2_____
lock(lock) {
    while(1) {
        old = ll(lock);
        if(old == 0)
            if(sc(lock, 1))
                return;
    }
}
```

# LLSC Lock Action Zone II



P1 State Data — lock: M  1

P2 State Data — lock: I

lock:  0

```
_____P1_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

```
_____P2_____
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

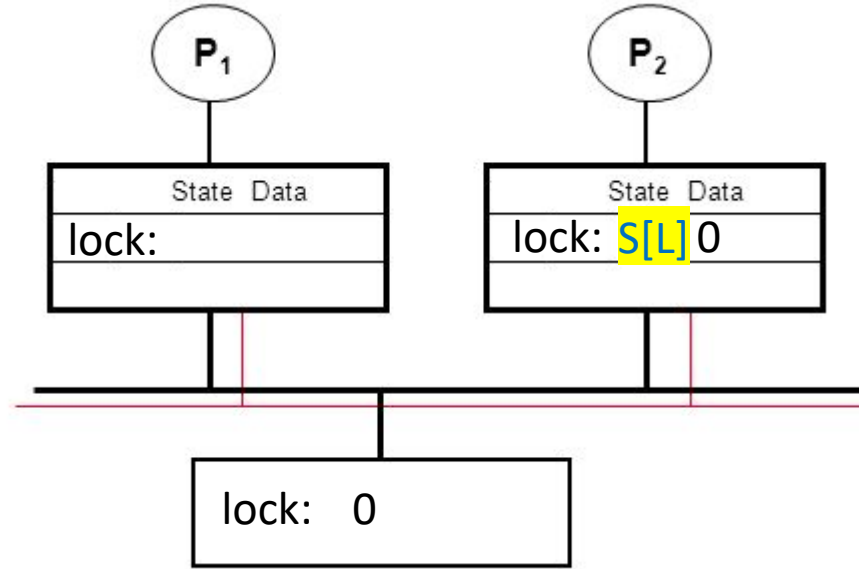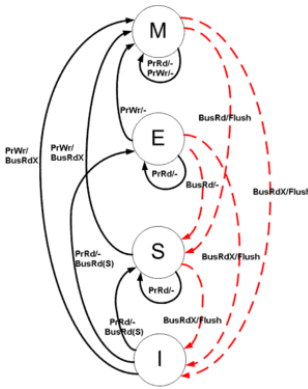# LLSC Lock Action Zone II



P1
```
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

P2
```
lock(lock) {
  while(1) {
    old = ll(lock);
    if(old == 0)
      if(sc(lock, 1))
        return;
  }
}
```

Store conditional fails

# Implementing Locks with Test&set

```
int lock_value = 0;
int* lock = &lock_value;
```

# Implementing Locks with Test&set

```
int lock_value = 0;
int* lock = &lock_value;
```

```
Lock::Acquire() {
while (test&set(lock) == 1)
  ; //spin
}
```

(test & set ~ CAS ~ LLSC)

# Implementing Locks with Test&set

```
int lock_value = 0;
int* lock = &lock_value;
```

```
Lock::Acquire() {
while (test&set(lock) == 1)
  ; //spin
}
```

(test & set  ~ CAS ~ LLSC)

```
Lock::Release() {
    *lock = 0;
}
```

# Implementing Locks with Test&set

```
int lock_value = 0;
int* lock = &lock_value;
```

```
Lock::Acquire() {
while (test&set(lock) == 1)
  ; //spin
}
```

(test & set  ~ CAS ~ LLSC)

```
Lock::Release() {
    *lock = 0;
}
```

## What is the problem with this?

A. CPU usage  B. Memory usage C. Lock::Acquire() latency
D. Memory bus usage E. Does not work

# Test & Set with Memory Hierarchies

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

CPU C
```
// in critical region
```

CPU A
```
while(test&set(lock));
// in critical region
```

CPU B
```
while(test&set(lock));
```

L1 …

L1 …

L2
```
lock: 1
…
```

L2 …

Main Memory
```
0xF0 lock: 1
0xF4 …
```

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

CPU C
// in critical region

CPU A
while(test&set(lock));
// in critical region

CPU B
while(test&set(lock));

L1
...

L1
...

L2
lock: 1
...

L2
...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**CPU C**
```
// in critical region
```

**CPU A**
```
while(test&set(lock));
// in critical region
```

**CPU B**
```
while(test&set(lock));
```

L1 ...

L1 ...

L2
```
lock: 1
...
```

L2 ...

Main Memory
```
0xF0 lock: 1
0xF4 ...
```

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

What happens to lock variable's cache line when different CPUs contend?

CPU A
while(test&set(lock));
// in critical region

CPU B
while(test&set(lock));

L1
...

L1
...

L2
lock: 1
...

L2
...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Set with Memory Hierarchies

CPU C
// in critical region

Initially, lock held by CPU C

CPU A, B busy-waiting

What happens to lock variable's cache line when different CPUs contend?

CPU A
while(test&set(lock));
// in critical region

CPU B
while(test&set(lock));

L1

0xF0 lock: 1

L1

...

L2

lock: 1
...

L2

...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Set with Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

What happens to lock variable's cache line when different CPUs contend?

# Test & Set with Memory Hierarchies

**CPU C**
`// in critical region`

Initially, lock held by CPU C

CPU A, B busy-waiting

What happens to lock variable's cache line when different CPUs contend?

**CPU A**
`while(test&set(lock));`
`// in critical region`

**CPU B**
`while(test&set(lock));`

L1
> 0xF0 lock: 1

L1

L2
> lock: 1
> ...

L2
> ...

Main Memory
> 0xF0 lock: 1
> 0xF4 ...

# Test & Set with Memory Hierarchies

CPU C
`// in critical region`

Initially, lock held by CPU C

CPU A, B busy-waiting

What happens to lock variable's cache line when different CPUs contend?

CPU A
`while(test&set(lock));`
`// in critical region`

CPU B
`while(test&set(lock));`

L1

`0xF0 lock: 1`

L1

`...`

L2

`lock: 1`
`...`

L2

`...`

Main Memory

`0xF0 lock: 1`
`0xF4 ...`

- With bus-locking, lock prefix blocks *everyone*
- With CAS, LL-SC, cache line cache line *"ping pongs"* amongst contenders

# TTS: Reducing busy wait contention

### Test&Set

```
Lock::Acquire() {
while (test&set(lock) == 1);
}
```

Busy-wait on in-memory copy

```
Lock::Release() {
    *lock = 0;
}
```

### Test&Test&Set

```
Lock::Acquire() {
while(1) {
  while (*lock == 1) ; // spin just reading
  if (test&set(lock) == 0)  break;
}
}
```

Busy-wait on cached copy

```
Lock::Release() {
*lock = 0;
}
```

# TTS: Reducing busy wait contention

### Test&Set

```
Lock::Acquire() {
while (test&set(lock) == 1);
}
```

Busy-wait on in-memory copy

```
Lock::Release() {
    *lock = 0;
}
```

### Test&Test&Set

```
Lock::Acquire() {
while(1) {
    while (*lock == 1) ; // spin just reading
    if (test&set(lock) == 0)  break;
}
}
```
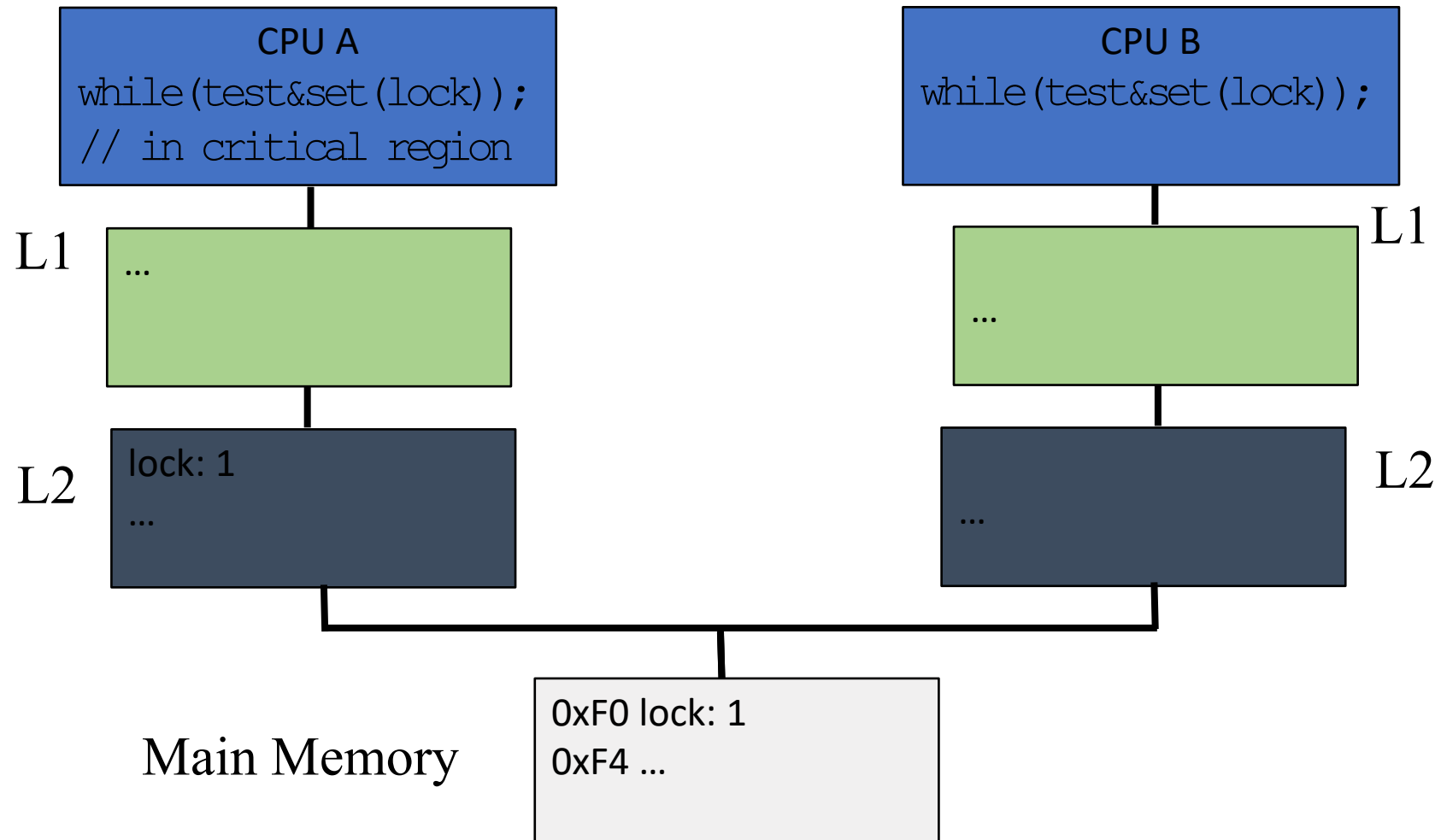
Busy-wait on cached copy

```
Lock::Release() {
*lock = 0;
}
```
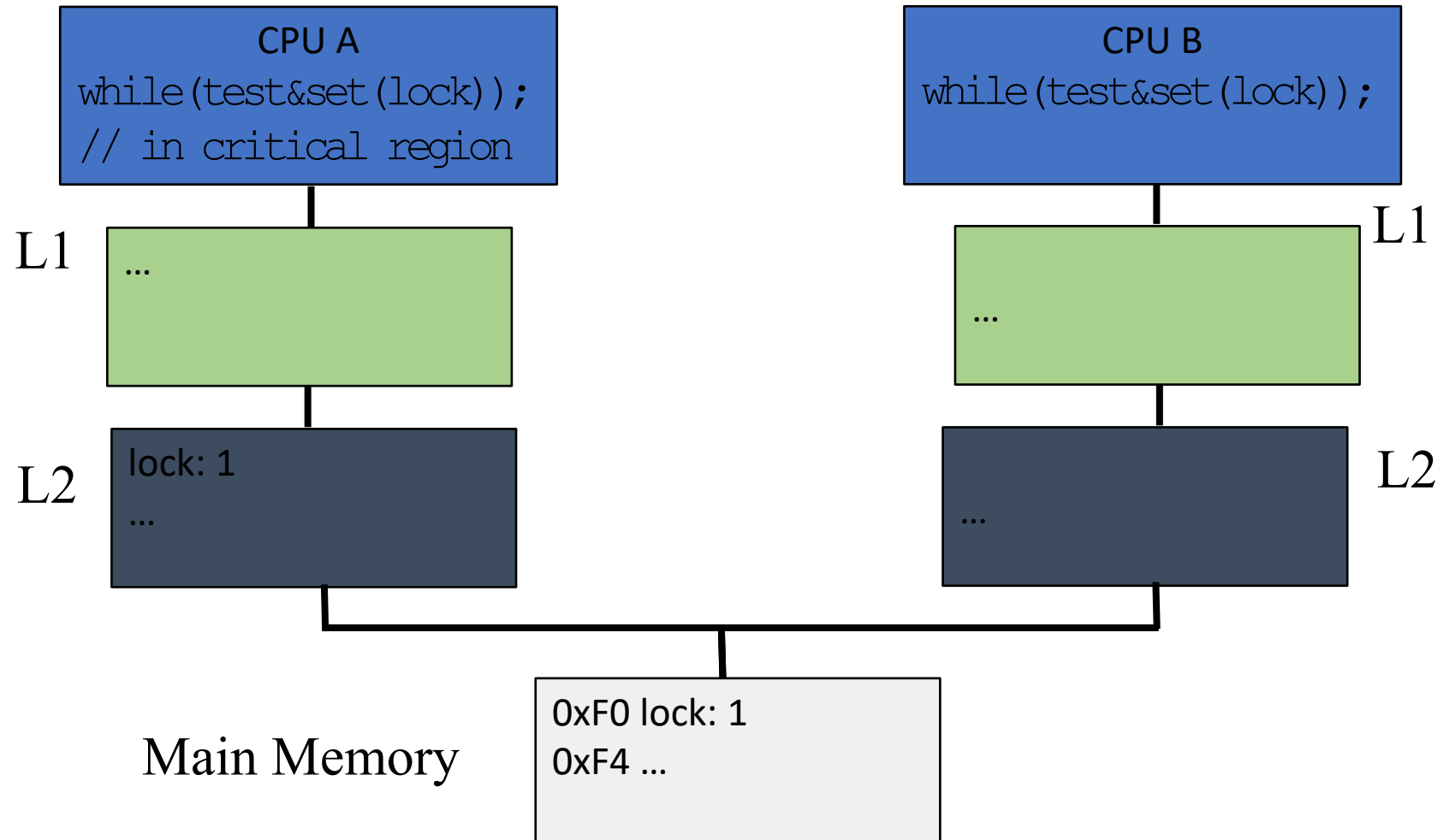
- What is the problem with this?
  - A. CPU usage  B. Memory usage C. Lock::Acquire() latency
  - D. Memory bus usage E. Does not work

# Test & Test & Set w Memory Hierarchies
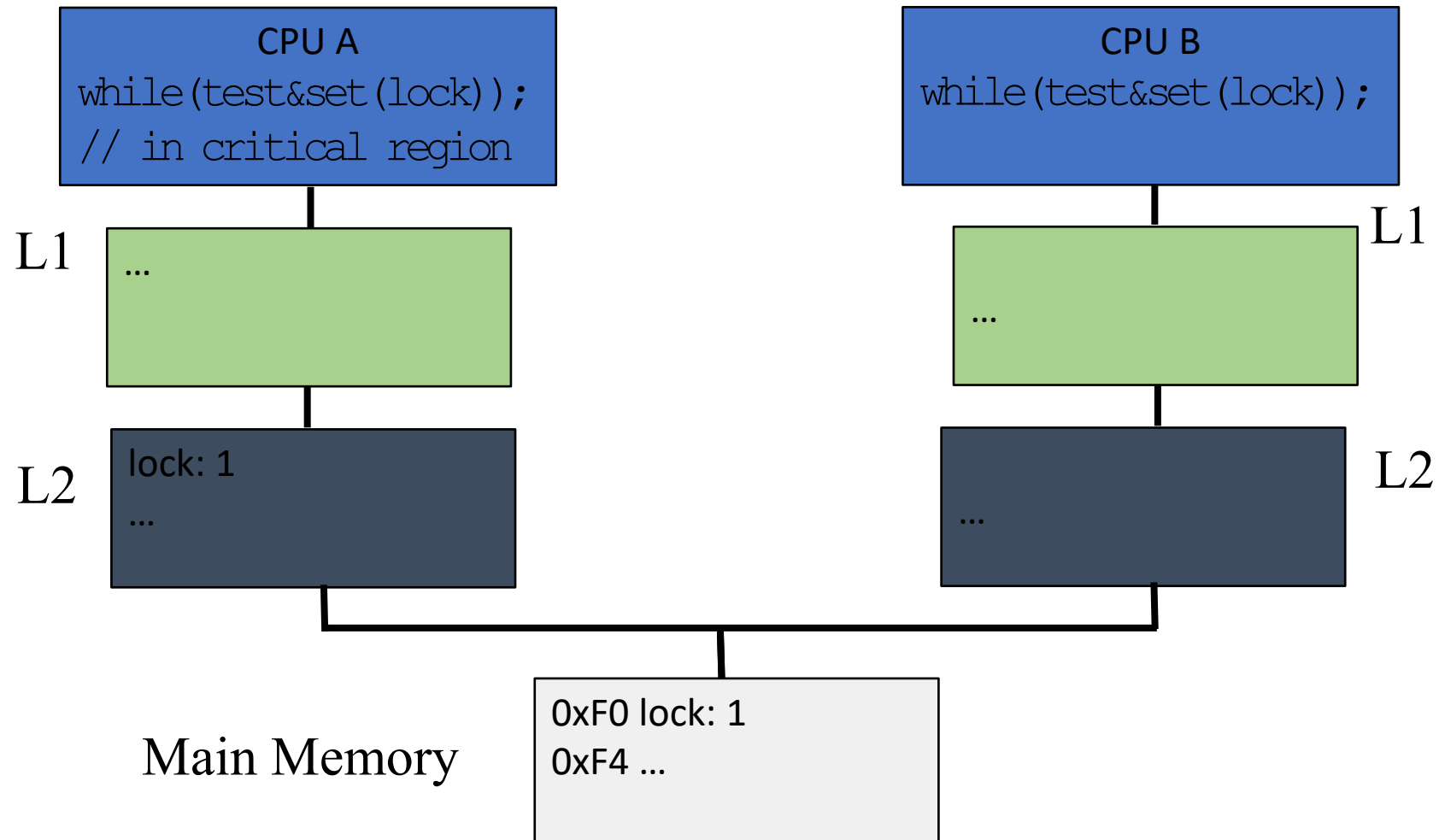
# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

# Test & Test & Set w Memory Hierarchies

CPU C
// in critical region

Initially, lock held by CPU C



CPU A
while(*lock);
if(test&set(lock))brk;

CPU B
while(*lock);
if(test&set(lock))brk;

L1
...

L1
...

L2
lock: 1
...

L2
...

Main Memory
0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1
...

L1
...

L2
lock: 1
...

L2
...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1
...

L1
...

L2
lock: 1
...

L2
...

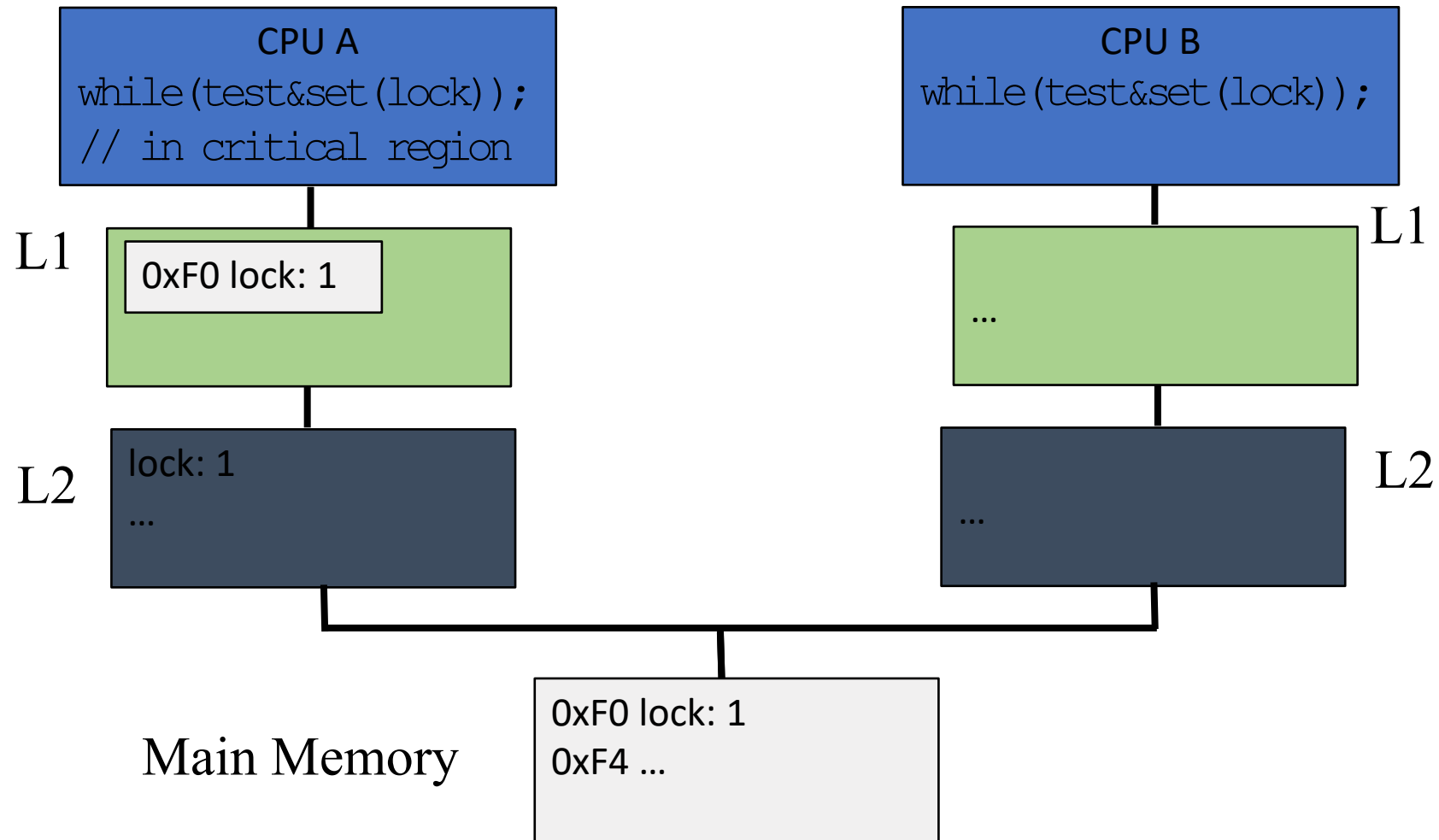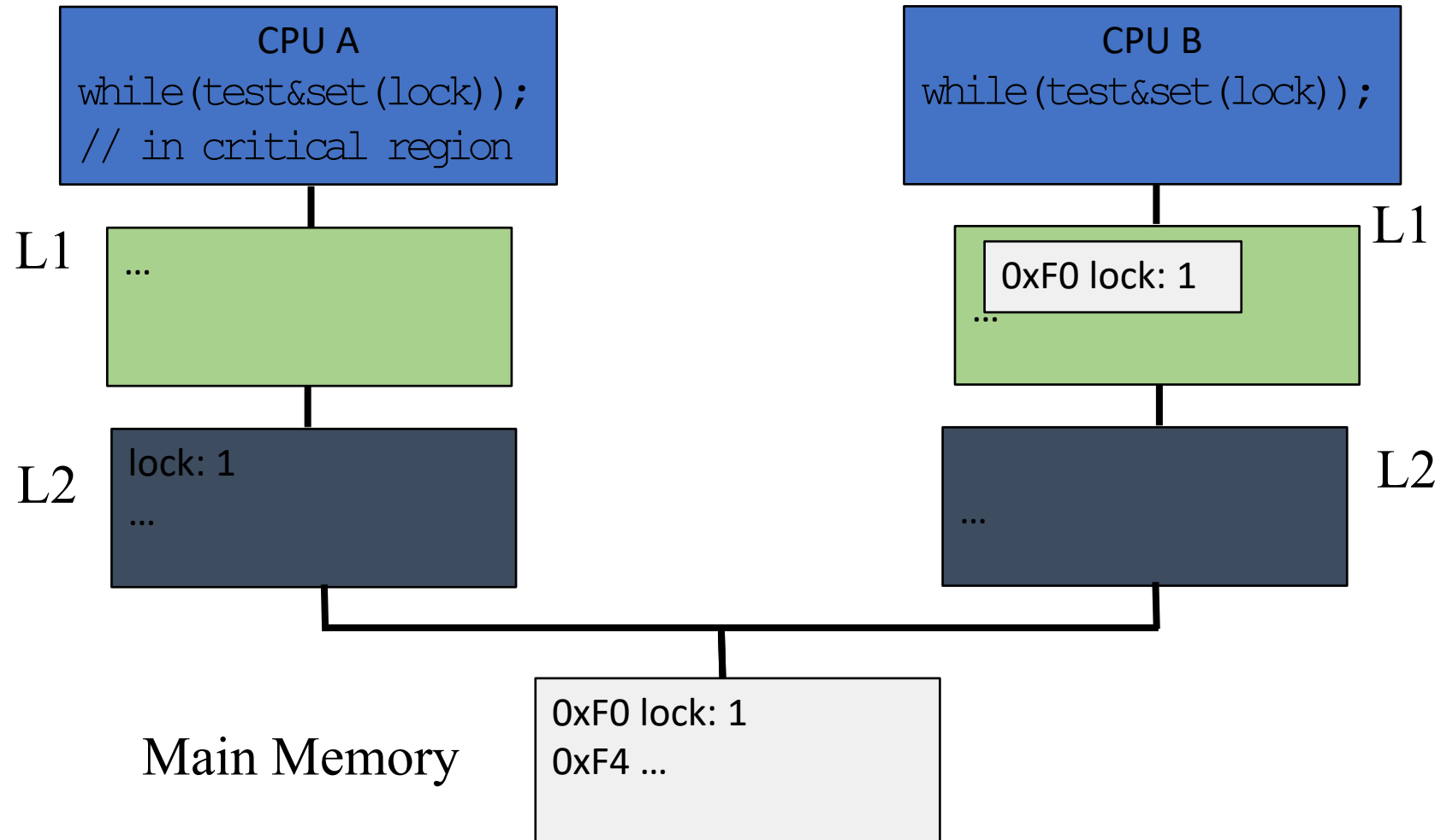Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

...

L1

...

L2

lock: 1
...

L2

...
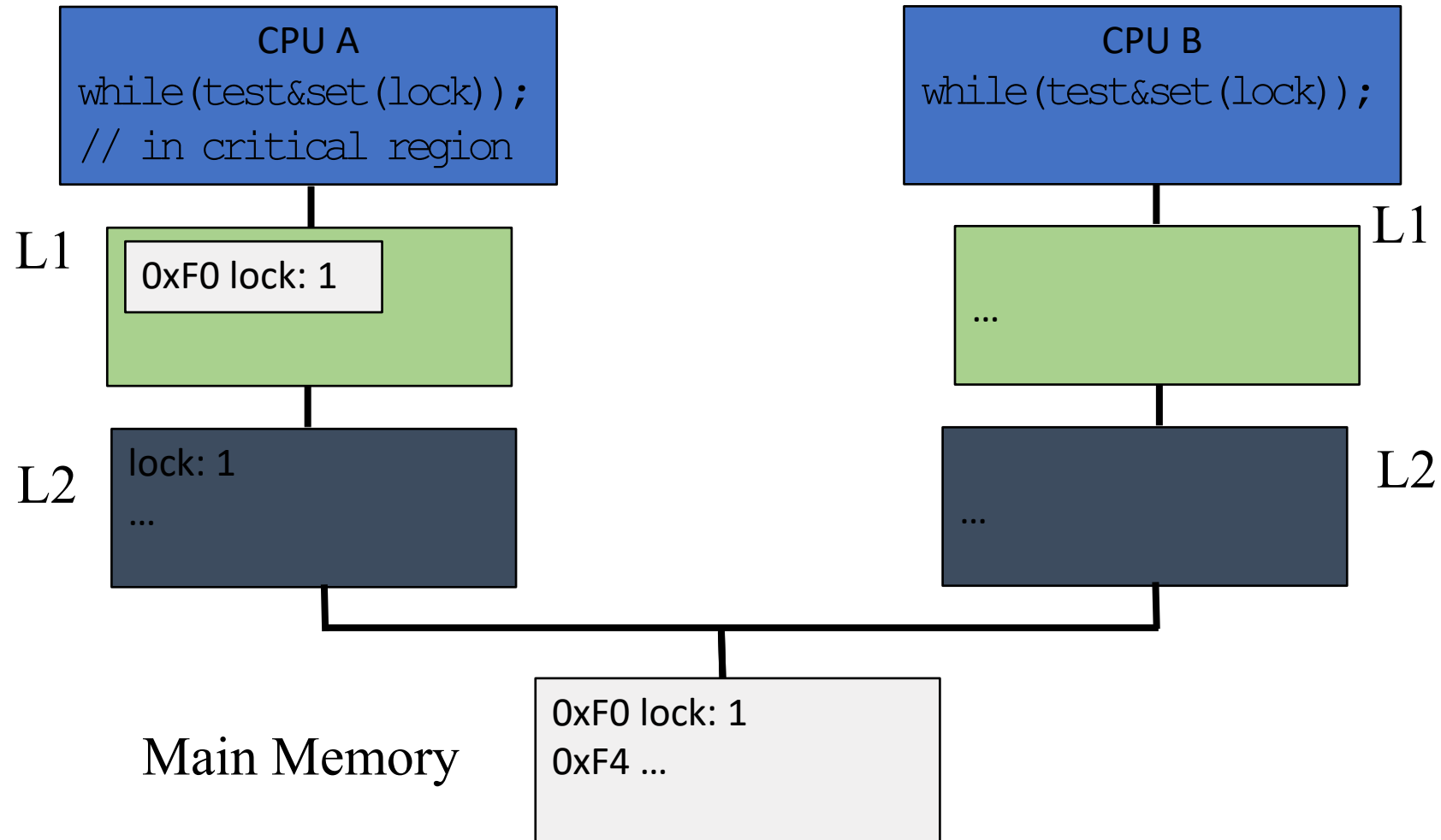
Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

0xF0 lock: 1

L1

...

L2

lock: 1
...

L2

...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies
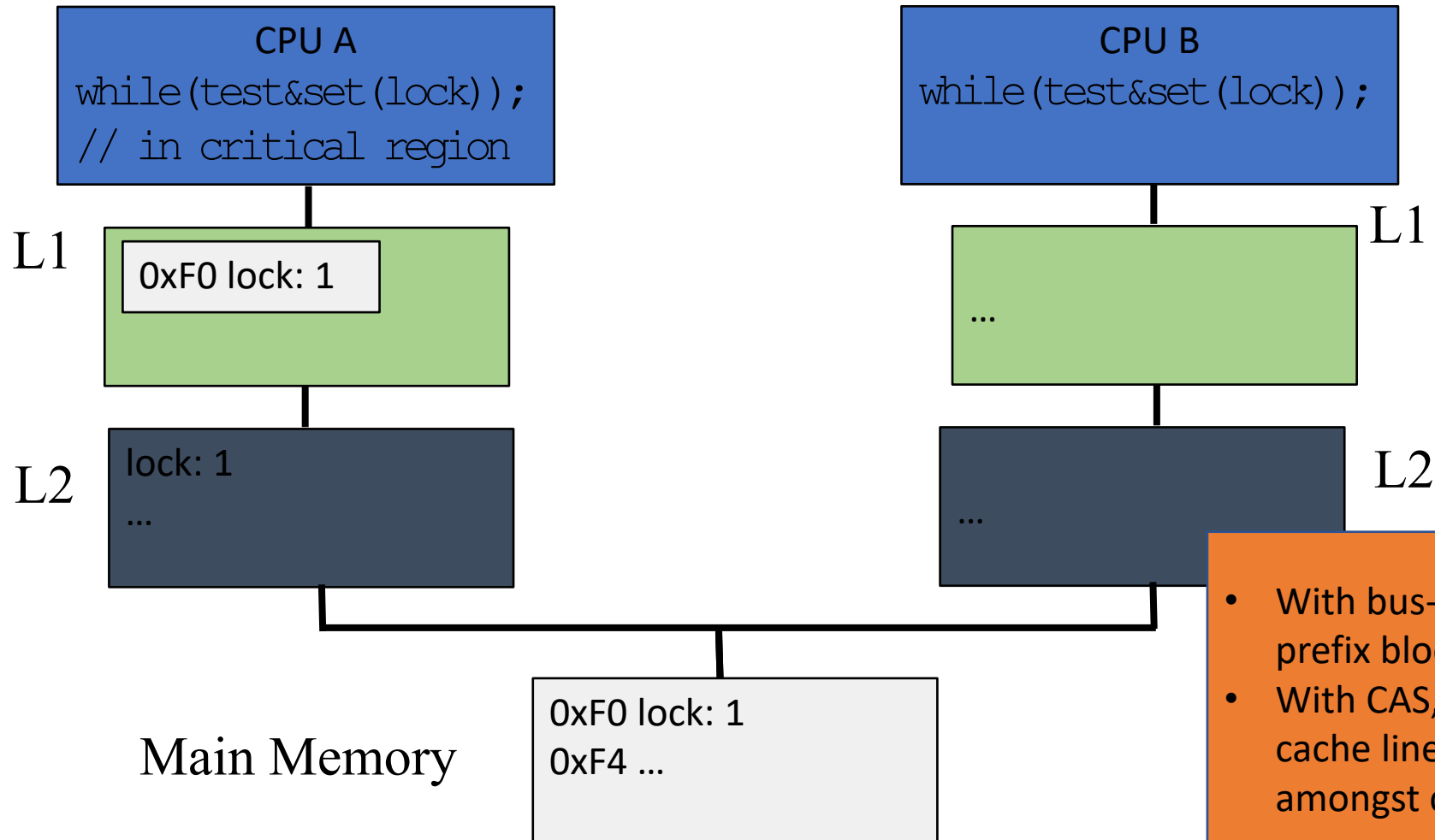
Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

0xF0 lock: 1

0xF0 lock: 1
...

L1

L2

lock: 1
...

...

L2

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?



CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

0xF0 lock: 1

L1

0xF0 lock: 1
...

**What coherence state is the lock in?**

L2

lock: 1
...

L2

...

0xF0 lock: 1
0xF4 ...

Main Memory

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

0xF0 lock: 1 **S**

L1

0xF0 lock: 1 **S**
...

**What coherence state is the lock in?**

L2

lock: 1
...

L2

...

Main Memory

0xF0 lock: 1
0xF4 ...

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

# Test & Test & Set w Memory Hierarchies

Initially, lock held by CPU C

CPU A, B busy-waiting

**Now** what happens to lock variable's cache line when different CPUs contend?

CPU A
```
while(*lock);
if(test&set(lock))brk;
```

CPU B
```
while(*lock);
if(test&set(lock))brk;
```

L1

0xF0 lock: 1 **S**

L1

0xF0 lock: 1 **S**
...

L2

lock: 1
...

L2

...

Main Memory

0xF0 lock: 1
0xF4 ...

- With TTS, all spin on own copy of cache line
- No more *"ping pongs"*

# How can we improve over busy-wait?

```
Lock::Acquire() {
while(1) {
  while (*lock == 1) ; // spin just reading
  if (test&set(lock) == 0)  break;
}
```

# Mutex

- Same abstraction as spinlock
- But is a "blocking" primitive
  - Lock available → same behavior
  - Lock held → yield/block
- Many ways to yield
- Simplest case of semaphore

```
void cm3_lock(u8_t* M) {
  u8_t LockedIn = 0;
  do {
    if (__LDREXB(Mutex) == 0) {
      // unlocked: try to obtain lock
      if ( __STREXB(1, Mutex)) { // got lock
        __CLREX(); // remove __LDREXB() lock
        LockedIn = 1;
      }
      else task_yield(); // give away cpu
    }
    else task_yield();    // give away cpu
} while(!LockedIn);
```

# Mutex

- Same abstraction as spinlock
- But is a "blocking" primitive
  - Lock available → same behavior
  - Lock held → yield/block
- Many ways to yield
- Simplest case of semaphore

```
void cm3_lock(u8_t* M) {
  u8_t LockedIn = 0;
  do {
    if (__LDREXB(Mutex) == 0) {
      // unlocked: try to obtain lock
      if ( __STREXB(1, Mutex)) { // got lock
        __CLREX(); // remove __LDREXB() lock
        LockedIn = 1;
      }
      else task_yield(); // give away cpu
    }
    else task_yield();   // give away cpu
  } while (!LockedIn);
```

- Is it better to use a spinlock or mutex on a uni-processor?

# Mutex

- Same abstraction as spinlock
- But is a "blocking" primitive
  - Lock available → same behavior
  - Lock held → yield/block
- Many ways to yield
- Simplest case of semaphore

```
void cm3_lock(u8_t* M) {
  u8_t LockedIn = 0;
  do {
    if (__LDREXB(Mutex) == 0) {
      // unlocked: try to obtain lock
      if ( __STREXB(1, Mutex)) { // got lock
        __CLREX(); // remove __LDREXB() lock
        LockedIn = 1;
      }
      else task_yield(); // give away cpu
    }
    else task_yield();   // give away cpu
} while (!LockedIn);
```

- Is it better to use a spinlock or mutex on a uni-processor?
- Is it better to use a spinlock or mutex on a multi-processor?

# Mutex

- Same abstraction as spinlock

- But is a "blocking" primitive
  - Lock available → same behavior
  - Lock held → yield/block

- Many ways to yield

- Simplest case of semaphore

```
void cm3_lock(u8_t* M) {
  u8_t LockedIn = 0;
  do {
    if (__LDREXB(Mutex) == 0) {
      // unlocked: try to obtain lock
      if ( __STREXB(1, Mutex)) { // got lock
        __CLREX(); // remove __LDREXB() lock
        LockedIn = 1;
      }
      else task_yield(); // give away cpu
    }
    else task_yield();    // give away cpu
} while (!LockedIn);
```

- Is it better to use a spinlock or mutex on a uni-processor?
- Is it better to use a spinlock or mutex on a multi-processor?
- How do you choose between spinlock/mutex on a multi-processor?

# Lock Pitfalls…

A(prio-0) → `lock(my_lock);`

B(prio-100) → `lock(my_lock);`

# Lock Pitfalls...

A(prio-0) →    `lock(my_lock);`

B(prio-100) → `lock(my_lock);`

**ACK! Priority Inversion!**

# Lock Pitfalls…



A(prio-0) →     `lock(my_lock);`

B(prio-100) → `lock(my_lock);`

**ACK! Priority Inversion!**

Solution?

# Lock Pitfalls...

A(prio-0) ➔ `lock(my_lock);`

B(prio-100) ➔ `lock(my_lock);`

<mark>**ACK! Priority Inversion!**</mark>

Solution?

**Priority inheritance:** A runs at B's priority

MARS pathfinder failure:

http://wiki.csie.ncku.edu.tw/embedded/priority-inversion-on-Mars.pdf

Other ideas?

# Can you build a lock without coherence?

# Can you build a lock without coherence?
# Dekker's Algorithm

```
variables
    wants_to_enter : array of 2 booleans
    turn : integer

wants_to_enter[0] ← false
wants_to_enter[1] ← false
turn ← 0    // or 1
```

```
p0:
   wants_to_enter[0] ← true
   while wants_to_enter[1] {
      if turn ≠ 0 {
         wants_to_enter[0] ← false
         while turn ≠ 0 {
            // busy wait
         }
         wants_to_enter[0] ← true
      }
   }

   // critical section
   ...
   turn ← 1
   wants_to_enter[0] ← false
   // remainder section
```

```
p1:
   wants_to_enter[1] ← true
   while wants_to_enter[0] {
      if turn ≠ 1 {
         wants_to_enter[1] ← false
         while turn ≠ 1 {
            // busy wait
         }
         wants_to_enter[1] ← true
      }
   }

   // critical section
   ...
   turn ← 0
   wants_to_enter[1] ← false
   // remainder section
```



initially: c1, c2, turn = 1,1,1

process 1

process 2

Th. J. Dekker's Solution

# Producer-Consumer (Bounded-Buffer) Problem

- Bounded buffer: size 'N'
  - Access entry 0… N-1, then "wrap around" to 0 again

- Producer process writes data to buffer
  - Must not write more than 'N' items more than consumer "consumes"

- Consumer process reads data from buffer
  - Should not try to consume if there is no data

# Producer-Consumer (Bounded-Buffer) Problem

- Bounded buffer: size 'N'
  - Access entry 0… N-1, then "wrap around" to 0 again

- Producer process writes data to buffer
  - Must not write more than 'N' items more than consumer "consumes"

- Consumer process reads data from buffer
  - Should not try to consume if there is no data

# Producer-Consumer (Bounded-Buffer) Problem

- Bounded buffer: size 'N'
  - Access entry 0… N-1, then "wrap around" to 0 again

- Producer process writes data to buffer
  - Must not write more than 'N' items more than consumer "consumes"

- Consumer process reads data from buffer
  - Should not try to consume if there is no data

# OK, let's write some code for this
(using locks only)

```
object array[N]
void enqueue(object x);
object dequeue();
```

# OK, let's write some code for this
(using locks only)

- Bounded buffer: size 'N'
  - Access entry 0... N-1, then "wrap around" to 0 again
- Producer writes data
- Consumer reads data

```
object array[N]
void enqueue(object x);
object dequeue();
```

# OK, let's write some code for this
(using locks only)

object array[N]
void enqueue(object x);
object dequeue();

Producer          Consumer

# Semaphore Motivation

# Semaphore Motivation

- Problem with locks: mutual exclusion, but *no ordering*

# Semaphore Motivation

- Problem with locks: mutual exclusion, but *no ordering*

- Inefficient for producer-consumer (and lots of other things)
  - Producer: creates a resource
  - Consumer: uses a resource
  - bounded buffer between them
  - You need synchronization for correctness, *and...*
  - Scheduling order:
    - producer waits if buffer full, consumer waits if buffer empty

# Semaphores

- Synchronization variable
  - Integer value
    - Can't access value directly
    - Must initialize to some value
      - sem_init(sem_t *s, int pshared, unsigned int value)
  - Two operations
    - sem_wait, or down(), P()
    - sem_post, or up(), V()

# Semaphores

- Synchronization variable
  - Integer value
    - Can't access value directly
    - Must initialize to some value
      - sem_init(sem_t *s, int pshared, unsigned int value)
  - Two operations
    - sem_wait, or down(), P()
    - sem_post, or up(), V()

```
int sem_wait(sem_t *s) {
    wait until value of semaphore s
        is greater than 0
    decrement the value of
        semaphore s by 1
}
```

```
int sem_post(sem_t *s) {
    increment the value of
        semaphore s by 1
    if there are 1 or more
        threads waiting, wake 1
}
```

# Semaphores

- Synchronization variable
  - Integer value
    - Can't access value directly
    - Must initialize to some value
      - sem_init(sem_t *s, int pshared, unsigned int value)
  - Two operations
    - sem_wait, or down(), P()
    - sem_post, or up(), V()

```
function V(semaphore S, integer I):
    [S ← S + I]
function P(semaphore S, integer I):
    repeat:
        if S ≥ I:
            S ← S – I
            break ]
```

```
int sem_wait(sem_t *s) {
    wait until value of semaphore s
        is greater than 0
    decrement the value of
        semaphore s by 1
}
```

```
int sem_post(sem_t *s) {
    increment the value of
        semaphore s by 1
    if there are 1 or more
        threads waiting, wake 1
}
```

# Semaphore Uses

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

# Semaphore Uses

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?
    - Binary semaphore: X=1

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

# Semaphore Uses

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?
    - Binary semaphore: X=1
    - ( Counting semaphore: X>1 )

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

# Semaphore Uses

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?
    - Binary semaphore: X=1
    - ( Counting semaphore: X>1 )

- Scheduling order
  - One thread waits for another

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

# Semaphore Uses

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?
    - Binary semaphore: X=1
    - ( Counting semaphore: X>1 )

- Scheduling order
  - One thread waits for another

```
//thread 0
… // 1st half of computation
sem_post(s);
```

```
// thread 1

sem_wait(s);
… //2nd half of computation
```

# Semaphore Uses

- Mutual exclusion
  - Semaphore as mutex
  - What should initial value be?
    - Binary semaphore: X=1
    - ( Counting semaphore: X>1 )

- Scheduling order
  - One thread waits for another
  - What should initial value be?

```
// initialize to X
sem_init(s, 0, X)

sem_wait(s);
// critical section
sem_post(s);
```

```
//thread 0
... // 1st half of computation
sem_post(s);
```

```
// thread 1

sem_wait(s);
... //2nd half of computation
```

# Producer-Consumer with semaphores

- Two semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots

# Producer-Consumer with semaphores

- Two semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots

```
sem_init(&full, 0, 0);
sem_init(&empty, 0, N);
```

```
producer() {
    sem_wait(empty);
    … // fill a slot
    sem_post(full);
}
```

```
consumer() {
    sem_wait(full);
    … // empty a slot
    sem_post(empty);
}
```

# Producer-Consumer with semaphores

- Two semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots

Is this correct?

```
sem_init(&full, 0, 0);
sem_init(&empty, 0, N);
```

```
producer() {
    sem_wait(empty);
    … // fill a slot
    sem_post(full);
}
```

```
consumer() {
    sem_wait(full);
    … // empty a slot
    sem_post(empty);
}
```

# Producer-Consumer with semaphores

- Two semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots

```
sem_init(&full, 0, 0);
sem_init(&empty, 0, N);
```

```
producer() {
    sem_wait(empty);
    ... // fill a slot
    sem_post(full);
}
```

```
consumer() {
    sem_wait(full);
    ... // empty a slot
    sem_post(empty);
}
```

# Producer-Consumer with semaphores

- Two semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots

- Problem: mutual exclusion?

```
sem_init(&full, 0, 0);
sem_init(&empty, 0, N);
```

```
producer() {
    sem_wait(empty);
    … // fill a slot
    sem_post(full);
}
```

```
consumer() {
    sem_wait(full);
    … // empty a slot
    sem_post(empty);
}
```

# Producer-Consumer with semaphores

- Three semaphores
  - sem_t full; // # of filled slots
  - sem_t empty; // # of empty slots
  - sem_t mutex; // mutual exclusion

```
sem_init(&full, 0, 0);
sem_init(&empty, 0, N);
sem_init(&mutex, 0, 1);
```

```
producer() {
    sem_wait(empty);
    sem_wait(&mutex);
    … // fill a slot
    sem_post(&mutex);
    sem_post(full);
}
```

```
consumer() {
    sem_wait(full);
    sem_wait(&mutex);
    … // empty a slot
    sem_post(&mutex);
    sem_post(empty);
}
```

# Pthreads and Semaphores

- **Type:** `pthread_semaphore_t`

```
int pthread_semaphore_init(pthread_spinlock_t *lock);
int pthread_semaphore_destroy(pthread_spinlock_t *lock);
…
```

- `?????`

# Pthreads and Semaphores

# Pthreads and Semaphores

- No pthread_semaphore_t!

# Pthreads and Semaphores

- No pthread_semaphore_t!
- POSIX does define standard

# Pthreads and Semaphores

- No pthread_semaphore_t!
- POSIX does define standard
- #include <semaphore.h>

- int **sem_wait**(sem_t ***sem**)
  - P action
  - blocks until the semaphore count pointed to by sem is greater than zero and then atomically decrements the count

- int **sem_post**(sem_t ***sem**)
  - V action
  - Atomically increments the count of the semaphore pointed to by sem. If there are any threads blocked on the semaphore, one will be unblocked

- int **sem_init(sem_t *sem**, int *pshared*, unsigned int *value*)
  - Initialize the semaphore to a value
  - If pshared is 0 then, semphamore is shared between threads of the process
    - else shared between processes

# What is a monitor?

# What is a monitor?

- ❑ Monitor: one big lock for set of operations/ methods
- ❑ Language-level implementation of mutex

- Entry procedure: called from outside
- Internal procedure: called within monitor
- Wait within monitor releases lock

Many variants…

# Pthreads and conditions/monitors

- **Type** `pthread_cond_t`

```
int pthread_cond_init(pthread_cond_t *cond,
                      const pthread_condattr_t *attr);
int pthread_cond_destroy(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond,
                      pthread_mutex_t *mutex);
int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
```

# Pthreads and conditions/monitors

- Type `pthread_cond_t`

```
int pthread_cond_init(pthread_cond_t *cond,
                          const pthread_condattr_t *attr);
int pthread_cond_destroy(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond,
                          pthread_mutex_t *mutex);
int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
```

# Pthreads and conditions/monitors

- Type `pthread_cond_t`

```
int pthread_cond_init(pthread_cond_t *cond,
                      const pthread_condattr_t *attr);
int pthread_cond_destroy(pthread_cond_t *cond);
int pthread_cond_wait(pthread_co
                      pthread_mut
int pthread_cond_signal(pthread_c
int pthread_cond_broadcast(pthrea
```

Java:
`synchronized keyword`
`wait()/notify()/notifyAll()`

C#: Monitor class
`Enter()/Exit()/`
`Pulse()/PulseAll()`

# Does this code work?

# Does this code work?

```java
1  public class SynchronizedQueue<T> {
2
3      public void enqueue(T item) {
4          lock.lock();
5          try {
6              if(head == tail - 1)
7                  notFull.wait();
8              Q[head] = item;
9              if(++head == MAX_Q)
10                 head = 0;
11             notEmpty.signal();
12         } finally {
13             lock.unlock();
14         }
15     }
16
17     public T dequeue() {
18         T retval = null;
19         lock.lock();
20         try {
21             if(head == tail)
22                 notEmpty.wait();
23             retval = Q[tail];
24             if(++tail == MAX_Q)
25                 tail = 0;
26             notFull.signal();
27         } finally {
28             lock.unlock();
29         }
30     }
31 }
```

# Does this code work?

```java
public class SynchronizedQueue<T> {

    public void enqueue(T item) {
        lock.lock();
        try {
            if(head == tail - 1)
                notFull.wait();
            Q[head] = item;
            if(++head == MAX_Q)
                head = 0;
            notEmpty.signal();
        } finally {
            lock.unlock();
        }
    }

    public T dequeue() {
        T retval = null;
        lock.lock();
        try {
            if(head == tail)
                notEmpty.wait();
            retval = Q[tail];
            if(++tail == MAX_Q)
                tail = 0;
            notFull.signal();
        } finally {
            lock.unlock();
        }
    }
}
```

```java
private Lock lock = new ReentrantLock();
private Condition notEmpty = lock.newCondition();
private Condition notFull = lock.newCondition();
private int head = 0;
private int tail = 0;
private int size = MAX_Q;
private T[] Q = new T[size];
```

# Does this code work?

```java
1  public class SynchronizedQueue<T> {
2
3      public void enqueue(T item) {
4          lock.lock();
5          try {
6              if(head == tail - 1)
7                  notFull.wait();
8              Q[head] = item;
9              if(++head == MAX_Q)
10                 head = 0;
11             notEmpty.signal();
12         } finally {
13             lock.unlock();
14         }
15     }
16
17     public T dequeue() {
18         T retval = null;
19         lock.lock();
20         try {
21             if(head == tail)
22                 notEmpty.wait();
23             retval = Q[tail];
24             if(++tail == MAX_Q)
25                 tail = 0;
26             notFull.signal();
27         } finally {
28             lock.unlock();
29         }
30     }
31 }
```

```java
private Lock lock = new ReentrantLock();
private Condition notEmpty = lock.newCondition();
private Condition notFull = lock.newCondition();
private int head = 0;
private int tail = 0;
private int size = MAX_Q;
private T[] Q = new T[size];
```

- Uses "if" to check invariants.

# Does this code work?

```java
public class SynchronizedQueue<T> {

    public void enqueue(T item) {
        lock.lock();
        try {
            if(head == tail - 1)
                notFull.wait();
            Q[head] = item;
            if(++head == MAX_Q)
                head = 0;
            notEmpty.signal();
        } finally {
            lock.unlock();
        }
    }

    public T dequeue() {
        T retval = null;
        lock.lock();
        try {
            if(head == tail)
                notEmpty.wait();
            retval = Q[tail];
            if(++tail == MAX_Q)
                tail = 0;
            notFull.signal();
        } finally {
            lock.unlock();
        }
    }
}
```

```java
private Lock lock = new ReentrantLock();
private Condition notEmpty = lock.newCondition();
private Condition notFull = lock.newCondition();
private int head = 0;
private int tail = 0;
private int size = MAX_Q;
private T[] Q = new T[size];
```

- Uses "if" to check invariants.
- Why doesn't **if** work?

# Does this code work?

```java
public class SynchronizedQueue<T> {

    public void enqueue(T item) {
        lock.lock();
        try {
            if(head == tail - 1)
                notFull.wait();
            Q[head] = item;
            if(++head == MAX_Q)
                head = 0;
            notEmpty.signal();
        } finally {
            lock.unlock();
        }
    }

    public T dequeue() {
        T retval = null;
        lock.lock();
        try {
            if(head == tail)
                notEmpty.wait();
            retval = Q[tail];
            if(++tail == MAX_Q)
                tail = 0;
            notFull.signal();
        } finally {
            lock.unlock();
        }
    }
}
```

```java
private Lock lock = new ReentrantLock();
private Condition notEmpty = lock.newCondition();
private Condition notFull = lock.newCondition();
private int head = 0;
private int tail = 0;
private int size = MAX_Q;
private T[] Q = new T[size];
```

- Uses "if" to check invariants.
- Why doesn't **if** work?
- How could we MAKE it work?

# Hoare-style Monitors
## (aka blocking condition variables)

Given entrance queue 'e', signal queue 's', condition var 'c'

```
enter:
    if(locked):
        e.push_back(thread)
    else
        lock
```

```
schedule:
    if s.any()
        t ← s.pop_first()
        t.run
    else if e.any()
        t ← e.pop_first()
        t. run
    else
        unlock // monitor unoccupied
```

```
wait C:

  C.q.push_back(thread)
  schedule  // block this thread
```

```
signal C :

   if (C.q.any())

        t = C.q.pop_front() // t → "the signaled thread"
        s.push_back(t)
        t.run
        // block this thread
```

- Leave calls schedule
- Signaler must wait, but gets priority over threads on entrance queue
- How is this different from Mesa monitors?
- Is s queue necessary?

# Mesa-style monitors

(aka non-blocking condition variables)

```
enter:
    if locked:
        e.push_back(thread)
        block
    else
        lock
```

```
schedule:
    if e.any()
        t ← e.pop_front
        t. run
    else
        unlock
```
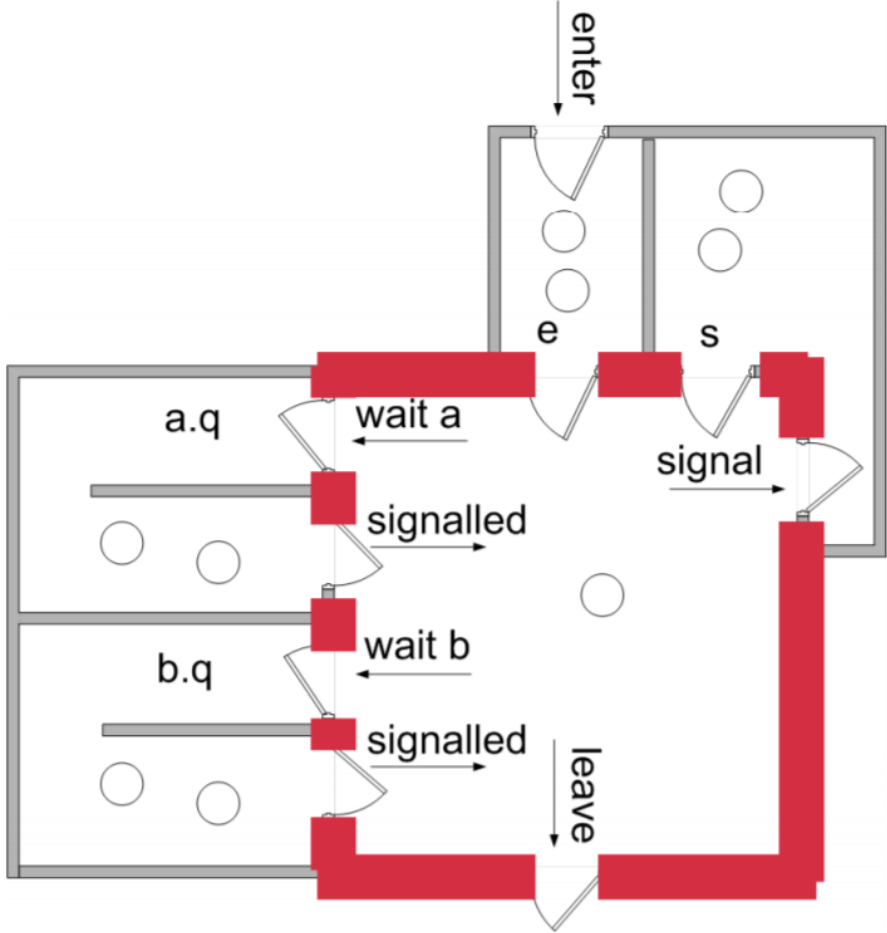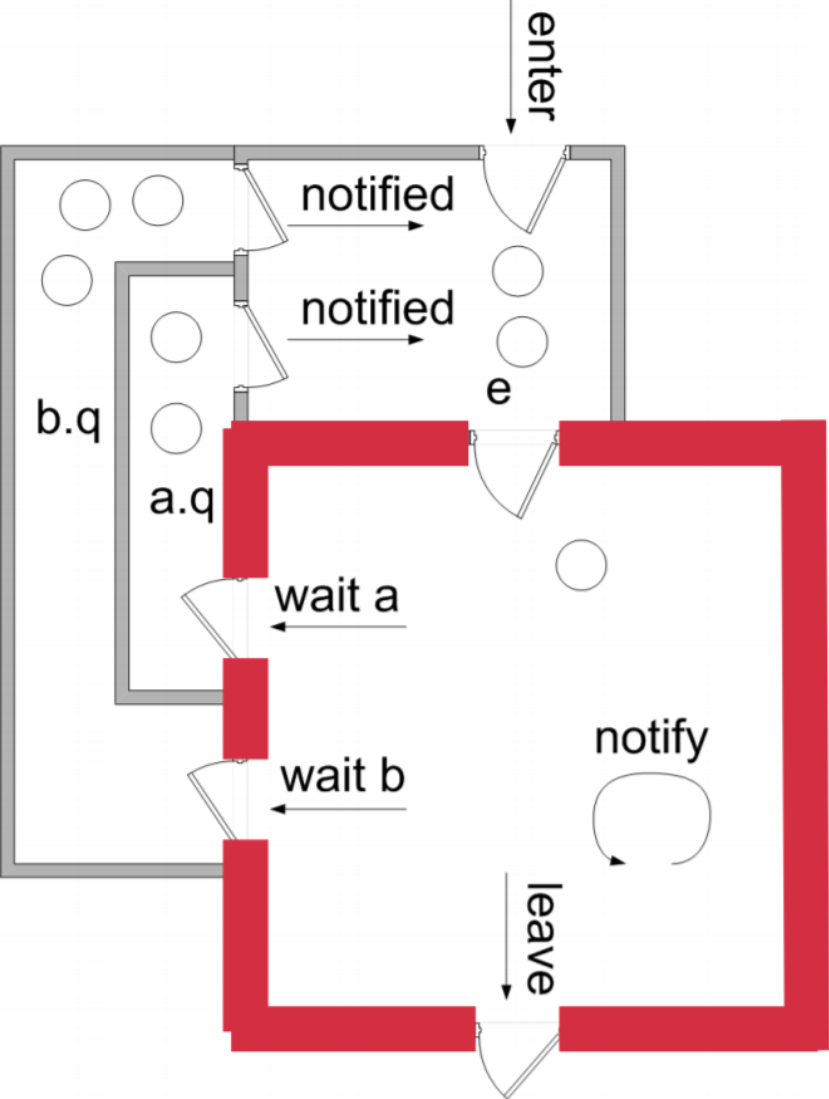
```
notify C:

    if C.q.any()

    t ← C.q.pop_front() // t is "notified "
    e.push_back(t)
```

```
wait C:

   C.q.push_back(thread)
   schedule
   block
```

- (Leave calls schedule)
- Can be extended with extra queues for priority
- What are the differences?

# Mesa, Hansen, Hoare

# Questions?