

CS/EE3810: Computer Organization

Lecture 17: Cache-Coherence and synchronization

Anton Burtsev
November, 2022

Memory organization

Memory Organization - I

- Centralized shared-memory multiprocessor or Symmetric shared-memory multiprocessor (SMP)
- Multiple processors connected to a single centralized memory – since all processors see the same memory organization ➡ uniform memory access (UMA)
- Shared-memory because all processors can access the entire memory address space
- Can centralized memory emerge as a bandwidth bottleneck? – not if you have large caches and employ fewer than a dozen processors

Cache Coherence Protocols

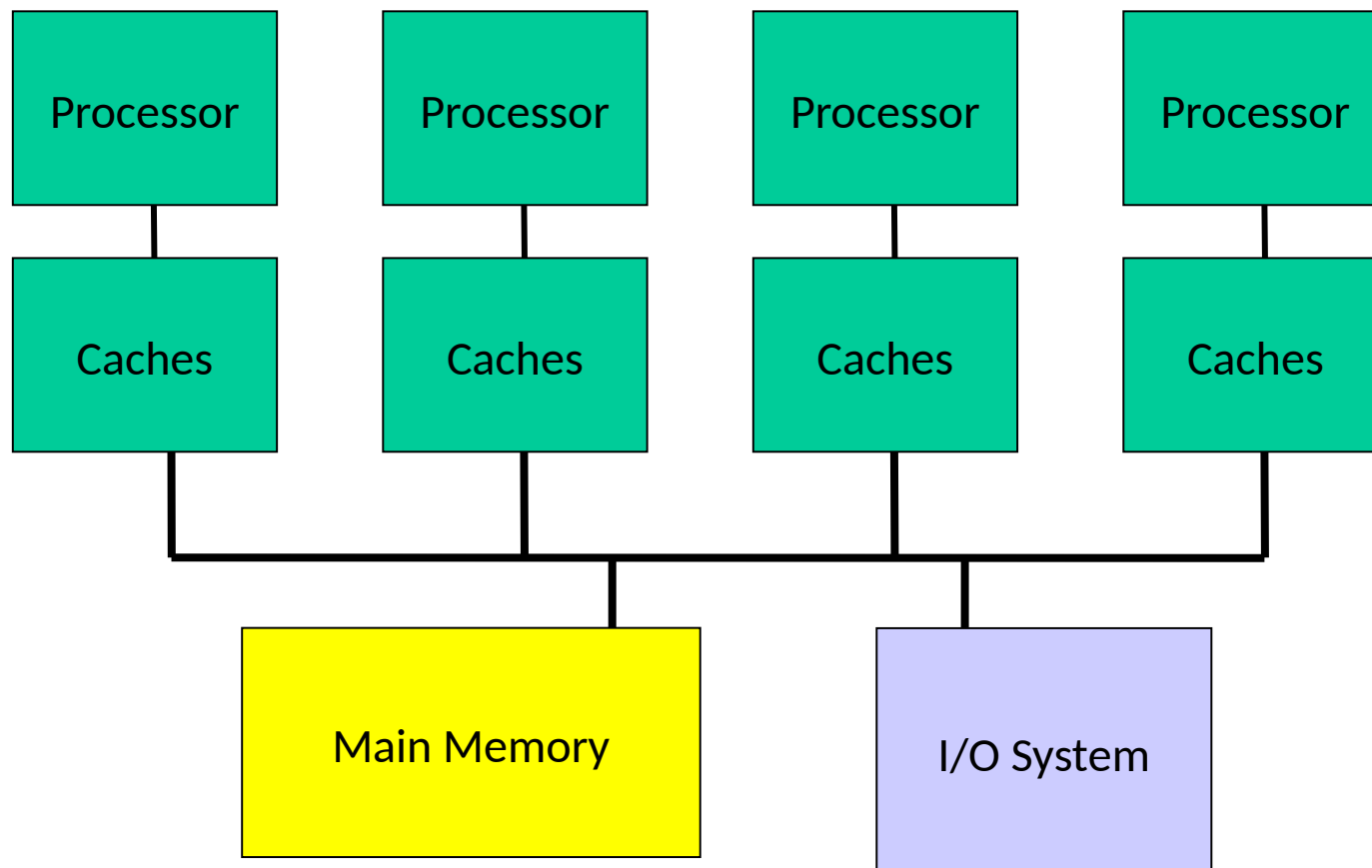
- Directory-based: A single location (directory) keeps track of the sharing status of a block of memory
- Snooping: Every cache block is accompanied by the sharing status of that block – all cache controllers monitor the shared bus so they can update the sharing status of the block, if necessary
- Write-invalidate: a processor gains exclusive access of a block before writing by invalidating all other copies
- Write-update: when a processor writes, it updates other shared copies of that block

Multiprocs -- Memory Organization - I

- Centralized shared-memory multiprocessor or Symmetric shared-memory multiprocessor (SMP)
- Multiple processors connected to a single centralized memory – since all processors see the same memory organization → uniform memory access (UMA)
- Shared-memory because all processors can access the entire memory address space
- Can centralized memory emerge as a bandwidth bottleneck? – not if you have large caches and employ fewer than a dozen processors

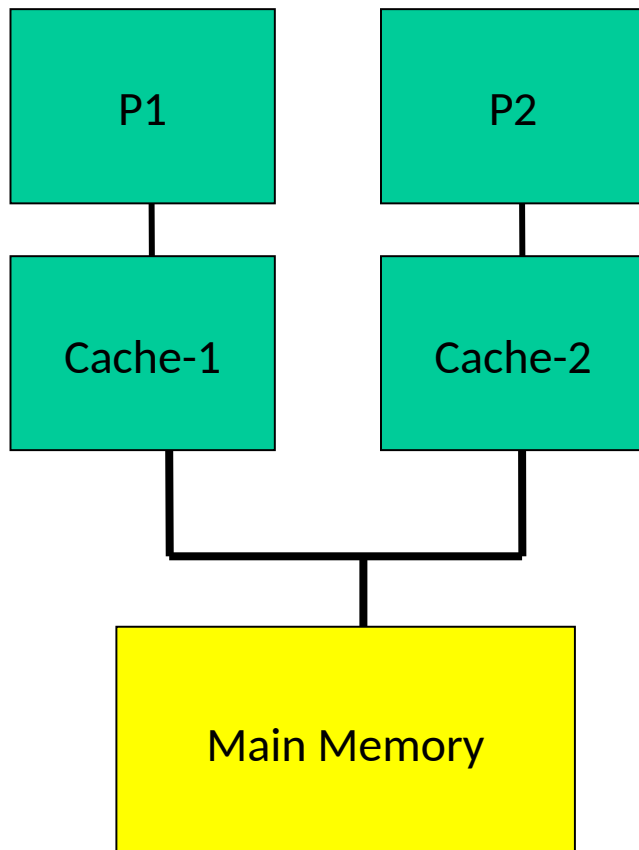
Snooping-Based Protocols

- Three states for a block: invalid, shared, modified
- A write is placed on the bus and sharers invalidate themselves
- The protocols are referred to as MSI, MESI, etc.



Example

- P1 reads X: not found in cache-1, request sent on bus, memory responds, X is placed in cache-1 in shared state
- P2 reads X: not found in cache-2, request sent on bus, everyone snoops this request, cache-1 does nothing because this is just a read request, memory responds, X is placed in cache-2 in shared state



- P1 writes X: cache-1 has data in shared state (shared only provides read perms), request sent on bus, cache-2 snoops and then invalidates its copy of X, cache-1 moves its state to modified
- P2 reads X: cache-2 has data in invalid state, request sent on bus, cache-1 snoops and realizes it has the only valid copy, so it downgrades itself to shared state and responds with data, X is placed in cache-2 in shared state, memory is also updated

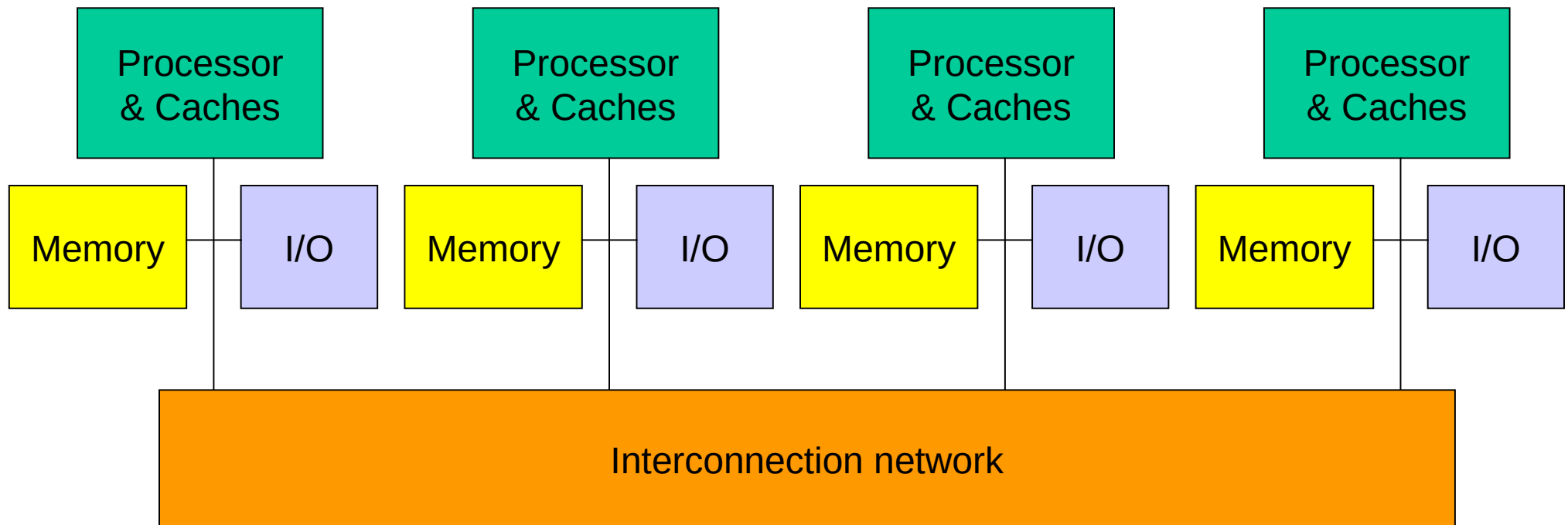
Example

Request	Cache Hit/Miss	Request on the bus	Who responds	State in Cache 1	State in Cache 2	State in Cache 3	State in Cache 4
				Inv	Inv	Inv	Inv
P1: Rd X	Miss	Rd X	Memory	S	Inv	Inv	Inv
P2: Rd X	Miss	Rd X	Memory	S	S	Inv	Inv
P2: Wr X	Perms Miss	Upgrade X	No response. Other caches invalidate.	Inv	M	Inv	Inv
P3: Wr X	Write Miss	Wr X	P2 responds	Inv	Inv	M	Inv
P3: Rd X	Read Hit	-	-	Inv	Inv	M	Inv
P4: Rd X	Read Miss	Rd X	P3 responds. Mem wrtbn	Inv	Inv	S	S

Multiprocs -- Memory Organization - II

- For higher scalability, memory is distributed among processors → distributed memory multiprocessors
- If one processor can directly address the memory local to another processor, the address space is shared → distributed shared-memory (DSM) multiprocessor
- If memories are strictly local, we need messages to communicate data → cluster of computers or multicomputers
- Non-uniform memory architecture (NUMA) since local memory has lower latency than remote memory

Distributed Memory Multiprocessors



Directory-based protocol

Request	Cache Hit/Miss	Messages	Dir State	State in C1	State in C2	State in C3	State in C4
				Inv	Inv	Inv	Inv
P1: Rd X	Miss	Rd-req to Dir. Dir responds.	X: S: 1	S	Inv	Inv	Inv
P2: Rd X	Miss	Rd-req to Dir. Dir responds.	X: S: 1, 2	S	S	Inv	Inv
P2: Wr X	Perms Miss	Upgr-req to Dir. Dir sends INV to P1. P1 sends ACK to Dir. Dir grants perms to P2.	X: M: 2	Inv	M	Inv	Inv
P3: Wr X	Write Miss	Wr-req to Dir. Dir fwds request to P2. P2 sends data to Dir. Dir sends data to P3.	X: M: 3	Inv	Inv	M	Inv
P3: Rd X	Read Hit	-	-	Inv	Inv	M	Inv
P4: Rd X	Read Miss	Rd-req to Dir. Dir fwds request to P3. P3 sends data to Dir. Memory wrtbk. Dir sends data to P4.	X: S: 3, 4	Inv	Inv	S	S

Synchronization

Constructing Locks

- Applications have phases (consisting of many instructions) that must be executed atomically, without other parallel processes modifying the data
- A lock surrounding the data/code ensures that only one program can be in a critical section at a time
- The hardware must provide some basic primitives that allow us to construct locks with different properties
- Lock algorithms assume an underlying cache coherence mechanism – when a process updates a lock, other processes will eventually see the update

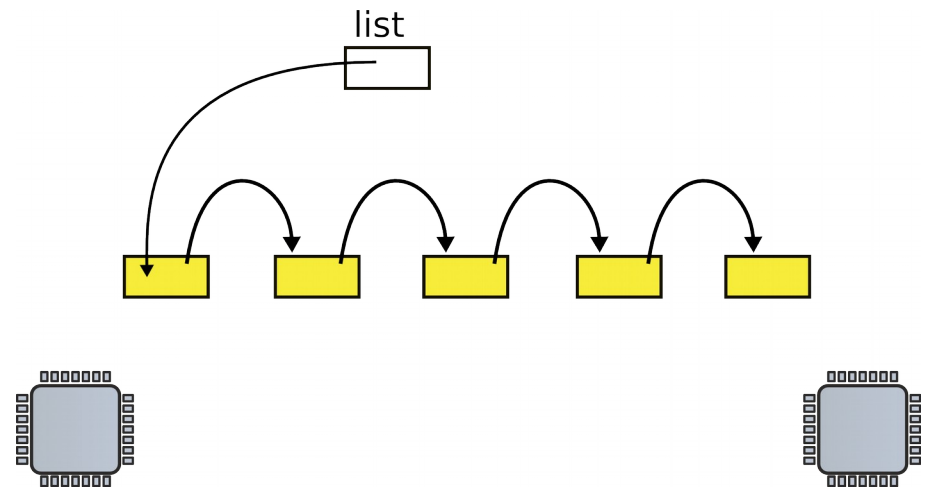
Race conditions

- Example:
 - Global list of, for example, requests
 - Each thread can add requests to the list

List implementation (no locks)

```
1 struct list {  
2     int data;  
3     struct list *next;  
4 };  
  
...  
6 struct list *list = 0;  
...  
9 insert(int data)  
10 {  
11     struct list *l;  
12  
13     l = malloc(sizeof *l);  
14     l->data = data;  
15     l->next = list;  
16     list = l;  
17 }
```

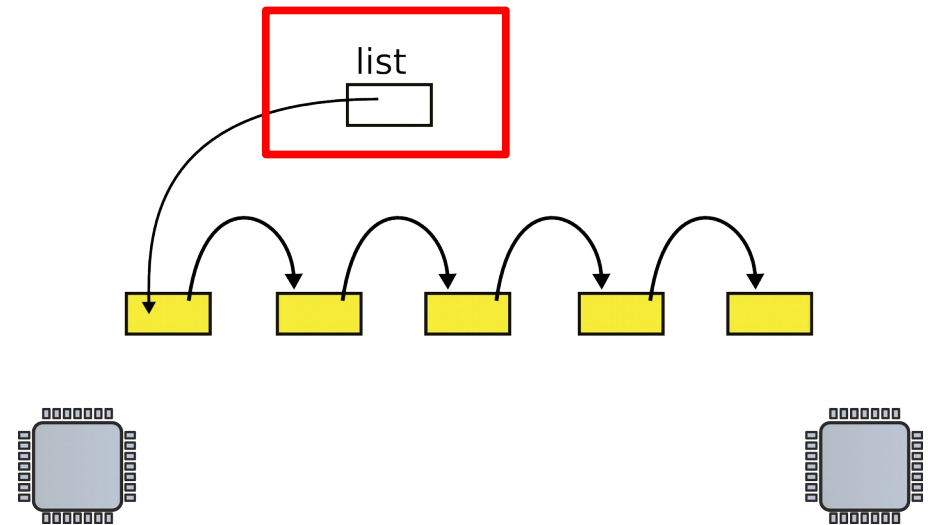
- List
 - One data element
 - Pointer to the next element



List implementation (no locks)

- Global head

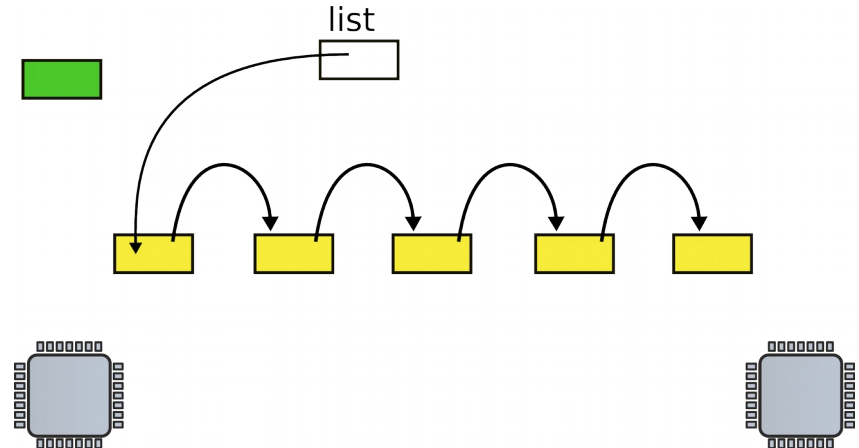
```
1 struct list {  
2     int data;  
3     struct list *next;  
4 };  
...  
6 struct list *list = 0;  
...  
9 insert(int data)  
10 {  
11     struct list *l;  
12  
13     l = malloc(sizeof *l);  
14     l->data = data;  
15     l->next = list;  
16     list = l;  
17 }
```



List implementation (no locks)

- Insertion
 - Allocate new list element

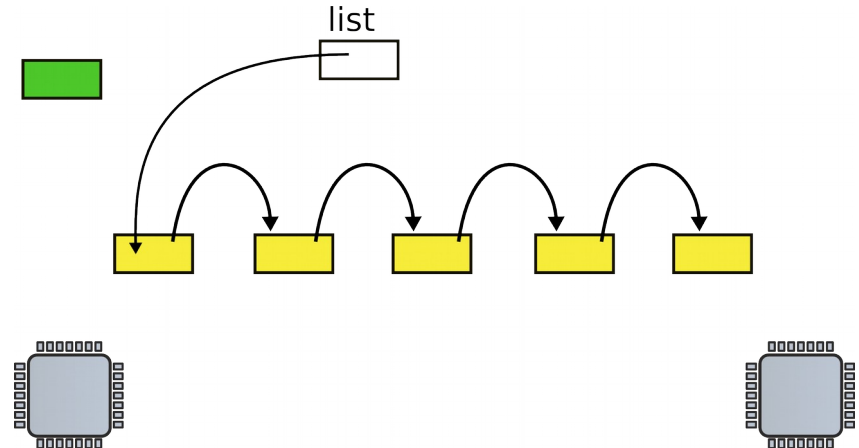
```
1 struct list {  
2     int data;  
3     struct list *next;  
4 };  
  
...  
6 struct list *list = 0;  
...  
9 insert(int data)  
10 {  
11     struct list *l;  
12  
13     l = malloc(sizeof *l);  
14     l->data = data;  
15     l->next = list;  
16     list = l;  
17 }
```



List implementation (no locks)

- Insertion
 - Allocate new list element
 - Save data into that element

```
1 struct list {  
2     int data;  
3     struct list *next;  
4 };  
  
...  
6 struct list *list = 0;  
...  
9 insert(int data)  
10 {  
11     struct list *l;  
12  
13     l = malloc(sizeof *l);  
14     l->data = data;  
15     l->next = list;  
16     list = l;  
17 }
```

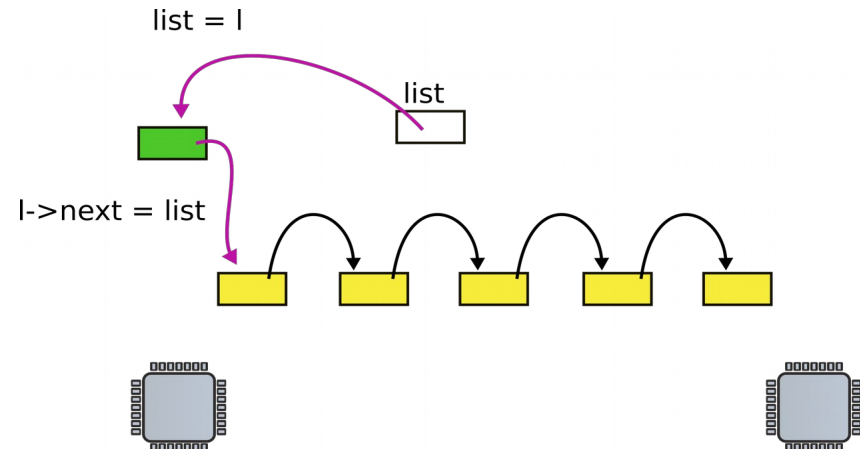


List implementation (no locks)

- Insertion

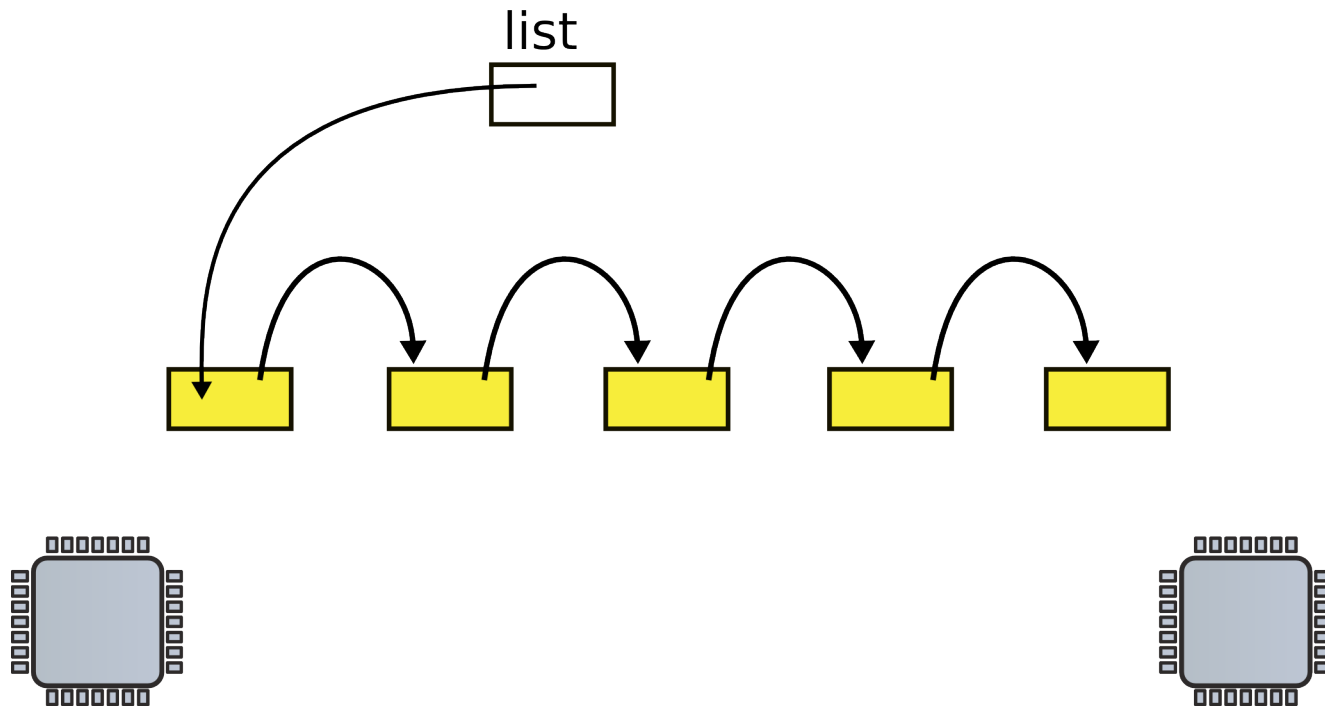
- Allocate new list element
- Save data into that element
- Insert into the list

```
1 struct list {
2     int data;
3     struct list *next;
4 };
...
6 struct list *list = 0;
...
9 insert(int data)
10 {
11     struct list *l;
12
13     l = malloc(sizeof *l);
14     l->data = data;
15     l->next = list;
16     list = l;
17 }
```



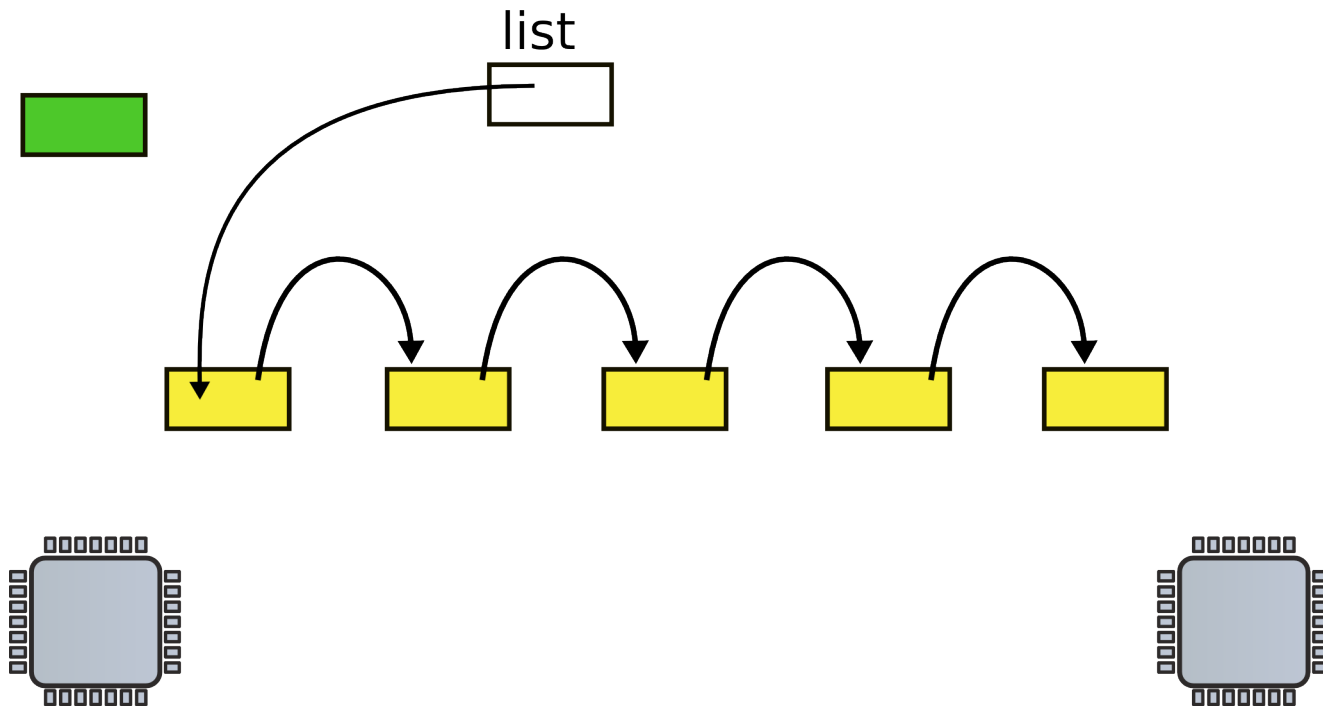
Now what happens when two CPUs access the
same list

Request queue (e.g. pending disk requests)

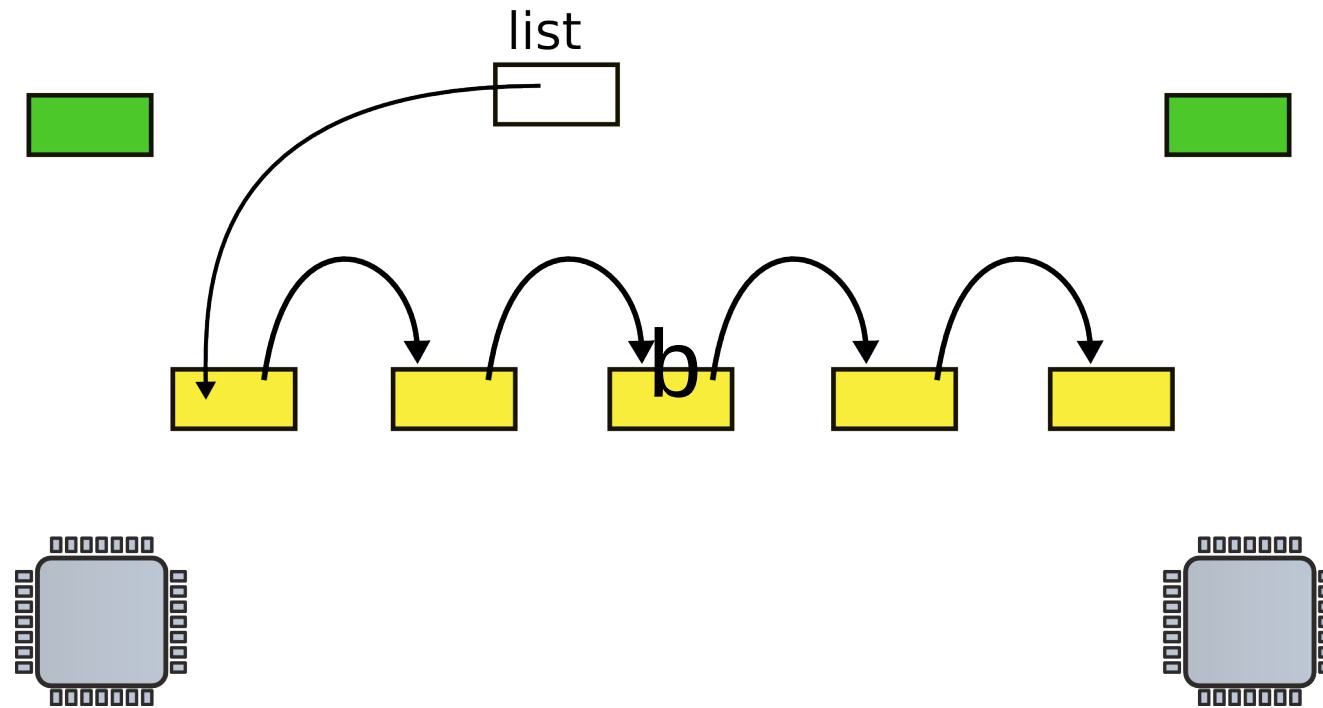


- Linked list, list is pointer to the first element

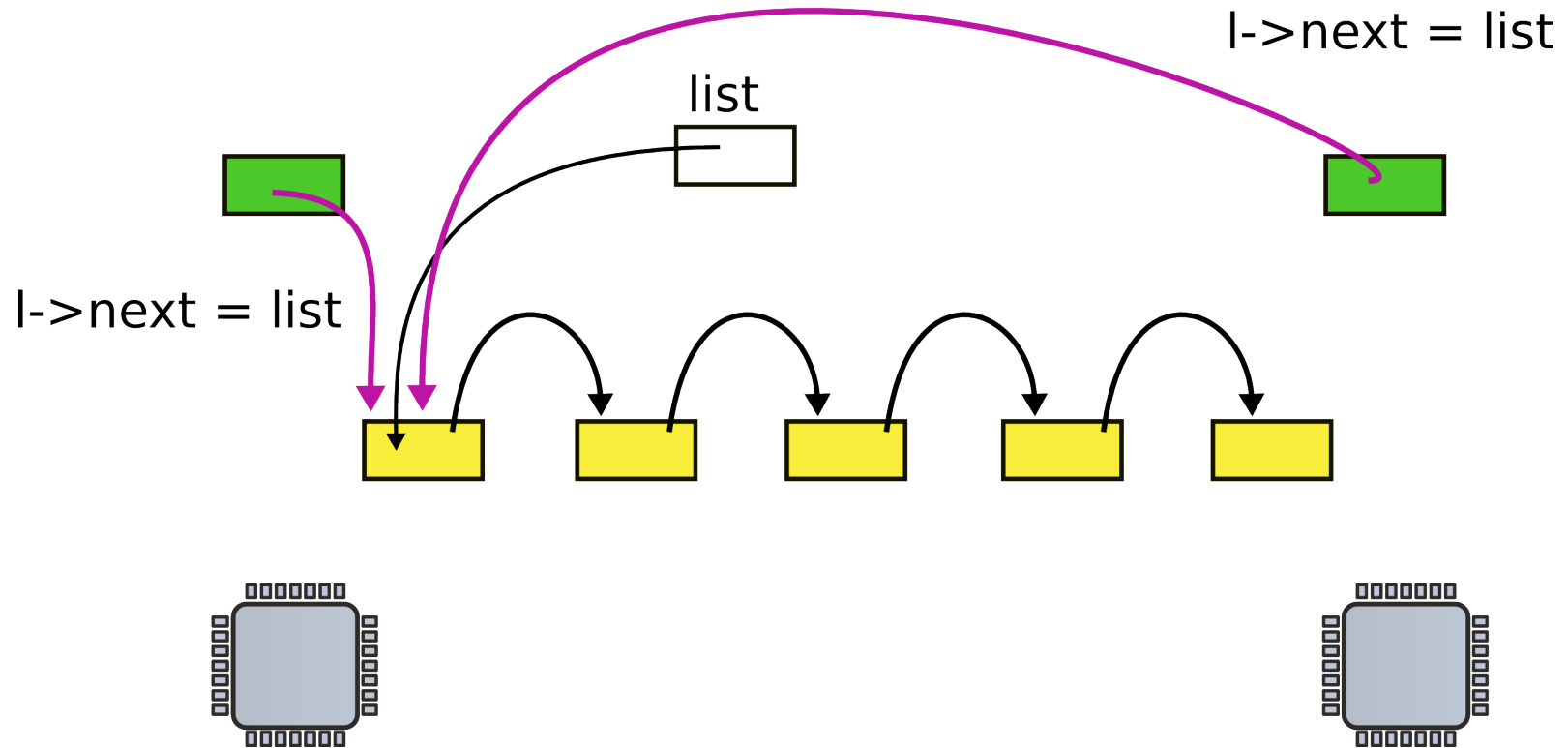
CPU1 allocates new request



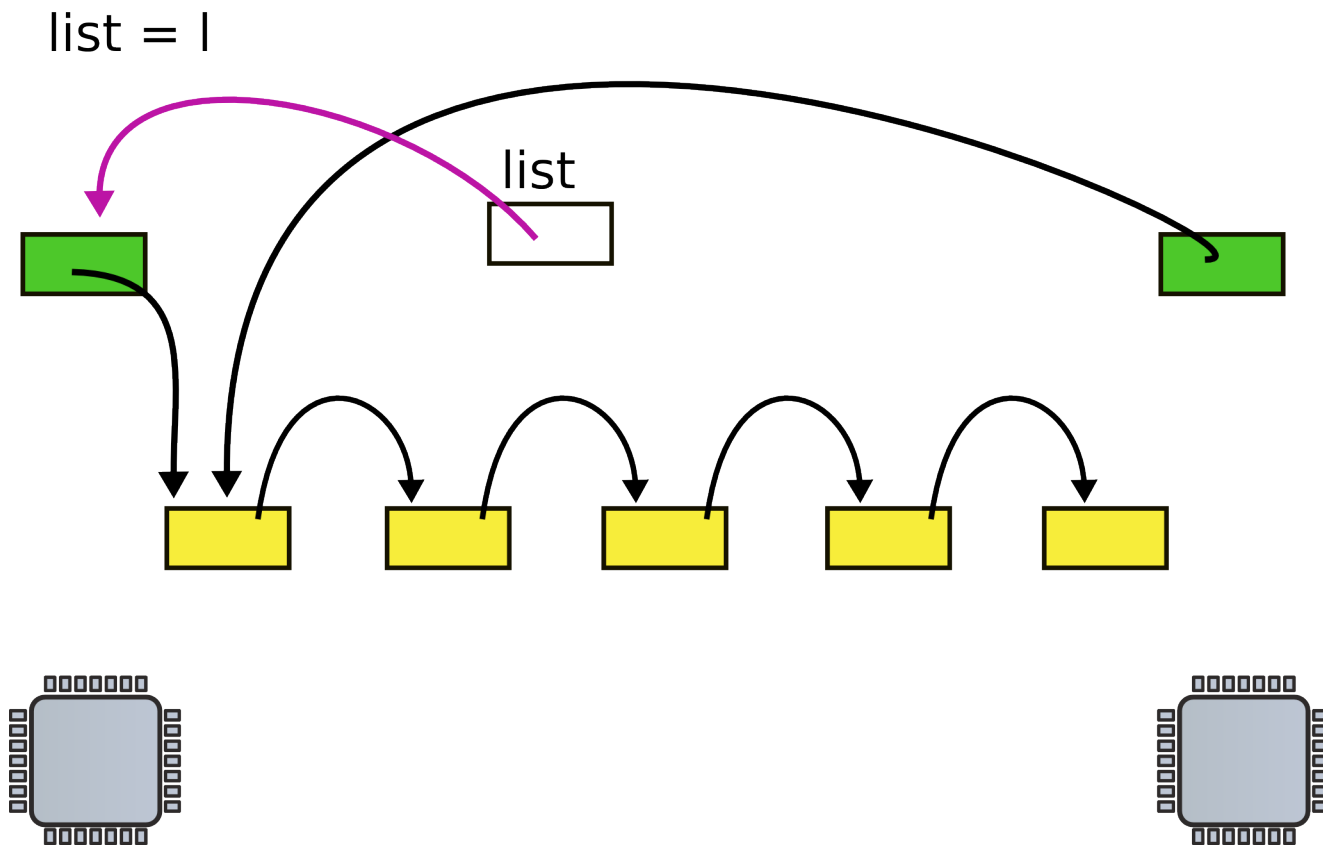
CPU2 allocates new request



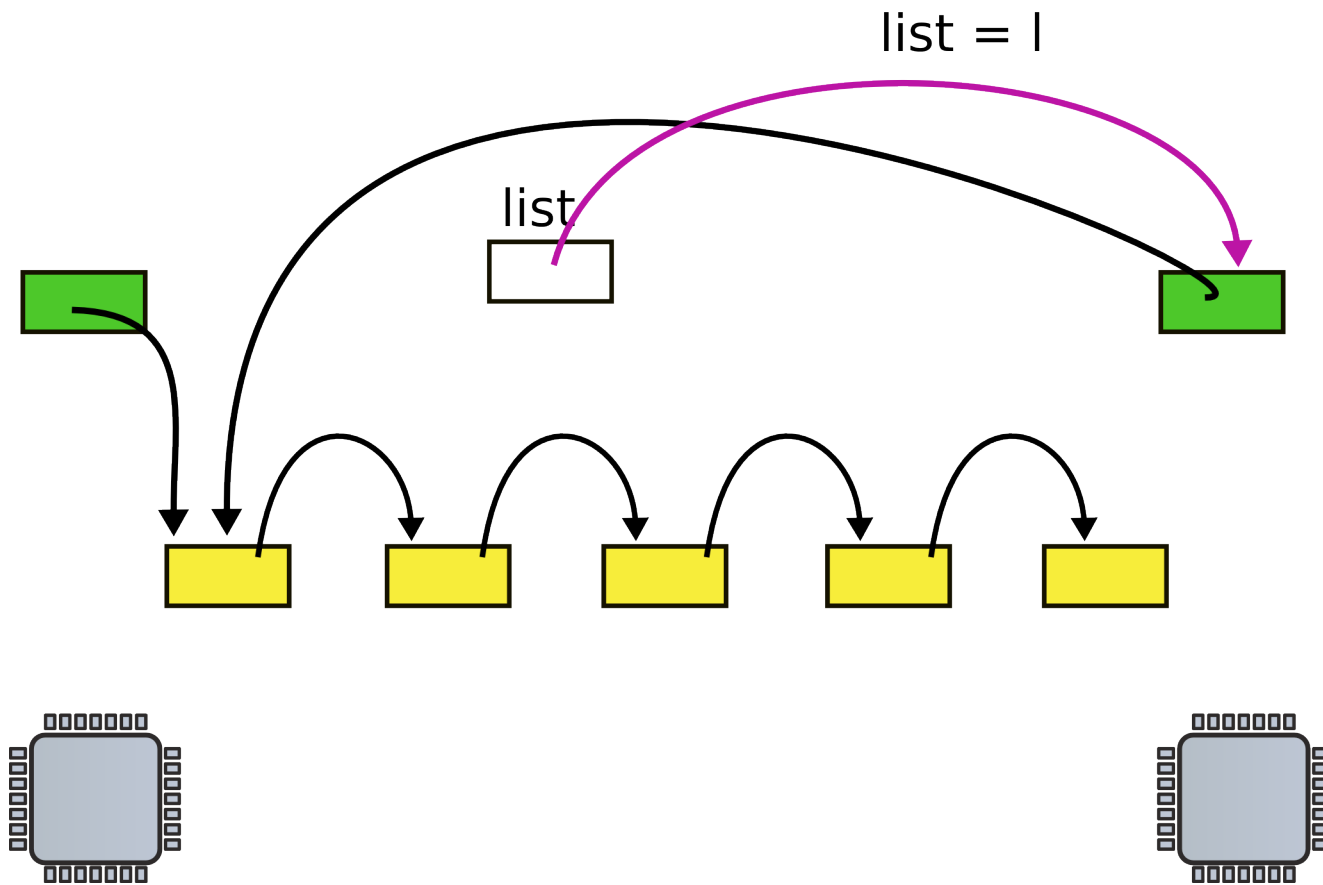
CPU 1 and 2 update next pointer



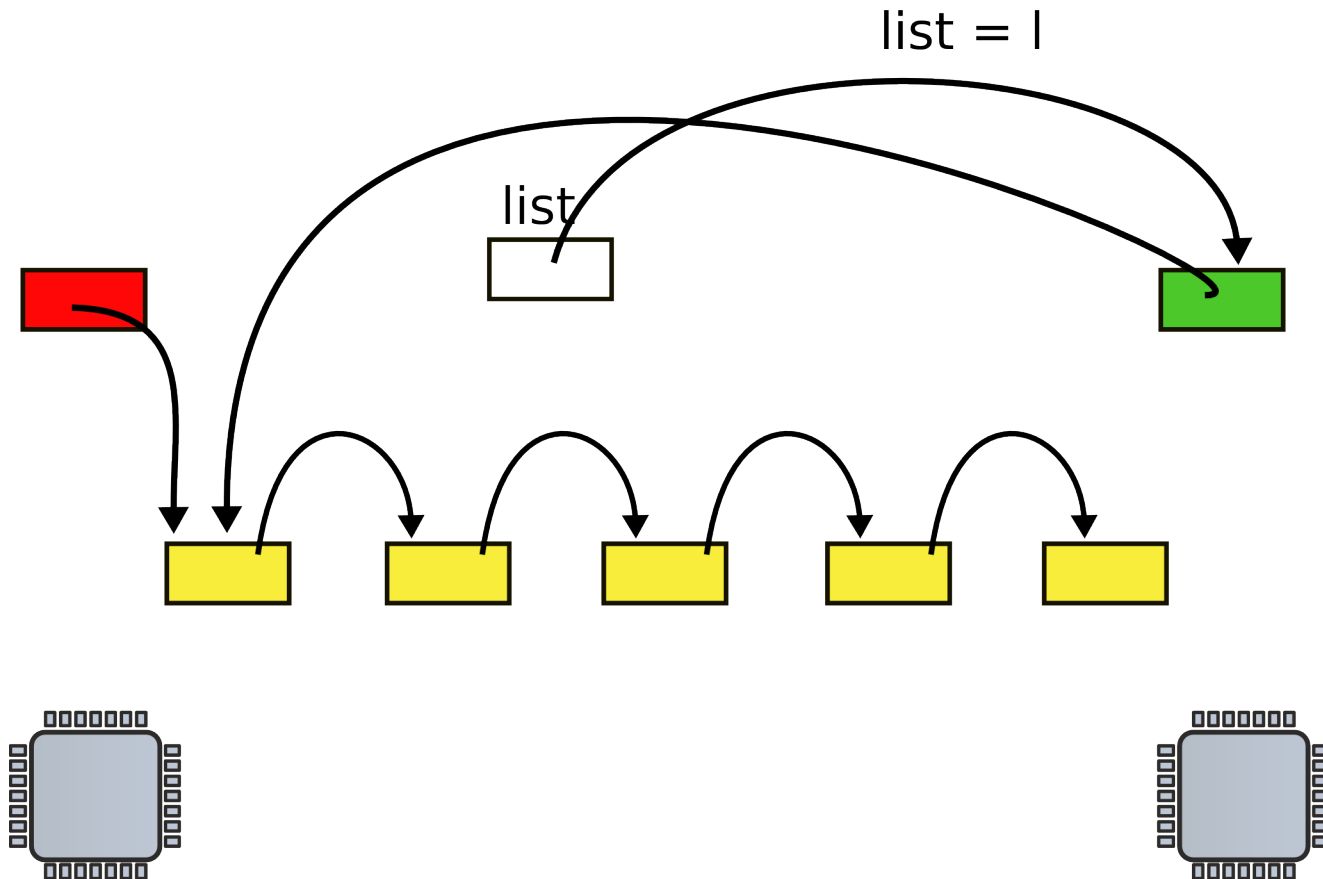
CPU1 updates head pointer



CPU2 updates head pointer



State after the race (red element is lost)



Mutual exclusion

- Only one CPU can update list at a time

List implementation with locks

```
1 struct list {
2     int data;
3     struct list *next;
4 };
5
6 struct list *list = 0;
7     struct lock listlock;
8
9 insert(int data)
10 {
11     struct list *l;
12
13     l = malloc(sizeof *l);
14     acquire(&listlock);
15
16     l->data = data;
17     l->next = list;
18     list = l;
19     release(&listlock);
20 }
21 }
```

- Critical section

- How can we implement `acquire()`?

Spinlock

```
21 void
22 acquire(struct spinlock *lk)
23 {
24     for(;;) {
25         if(!lk->locked) {
26             lk->locked = 1;
27             break;
28         }
29     }
30 }
```

- Spin until lock is 0
- Set it to 1

Still incorrect

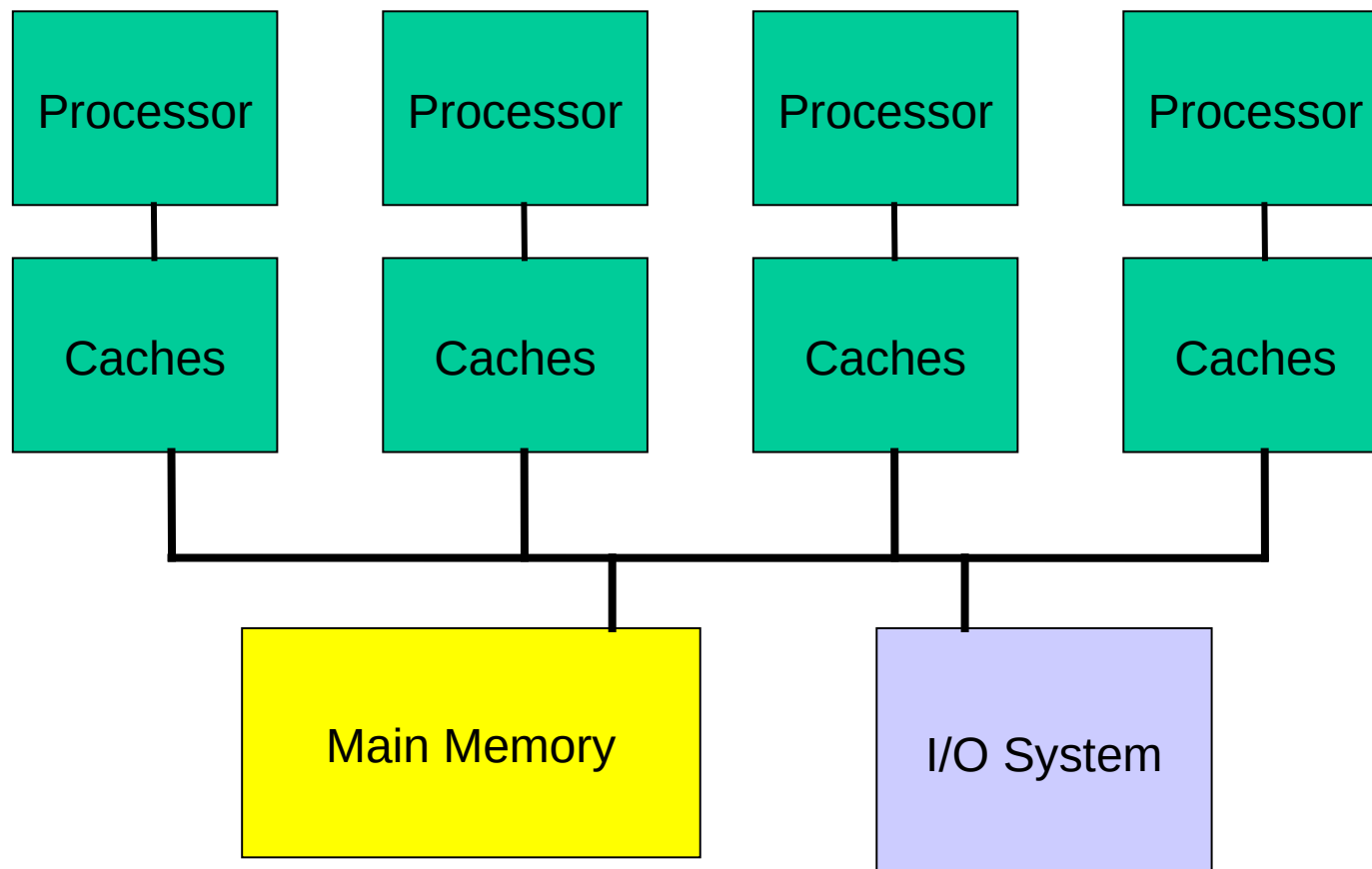
```
21 void
22 acquire(struct spinlock *lk)
23 {
24     for(;;) {
25         if(!lk->locked) {
26             lk->locked = 1;
27             break;
28         }
29     }
30 }
```

- Two CPUs can reach line #25 at the same time
 - See not locked, and
 - Acquire the lock
- Lines #25 and #26 need to be atomic
 - I.e. indivisible

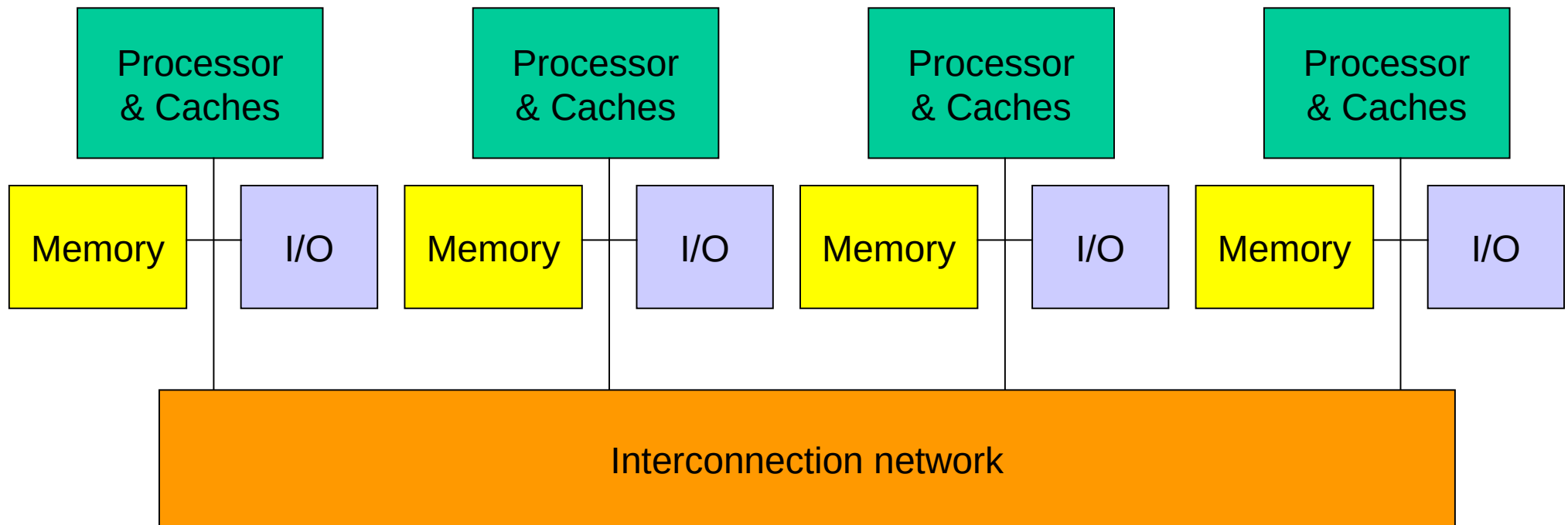
Synchronization

- The simplest hardware primitive that greatly facilitates synchronization implementations (locks, barriers, etc.) is an atomic read-modify-write
- Atomic exchange: swap contents of register and memory
- Special case of atomic exchange: test & set: transfer memory location into register and write 1 into memory
- acquire: t&s register, location
 bnz register, acquire
 CS
release: st location, #0

How does it work?



How does it work for directory based protocol

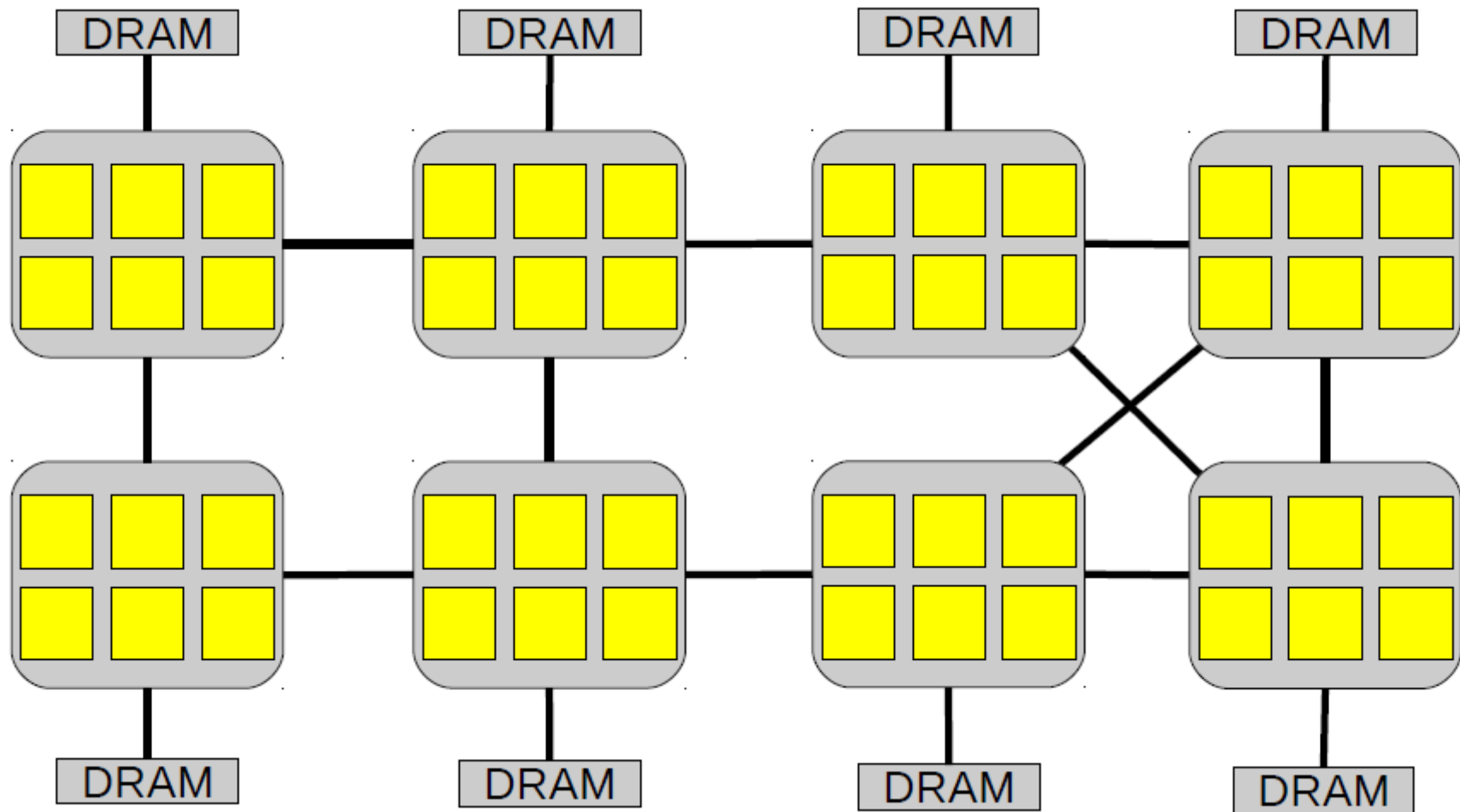


What is the main problem with locks?

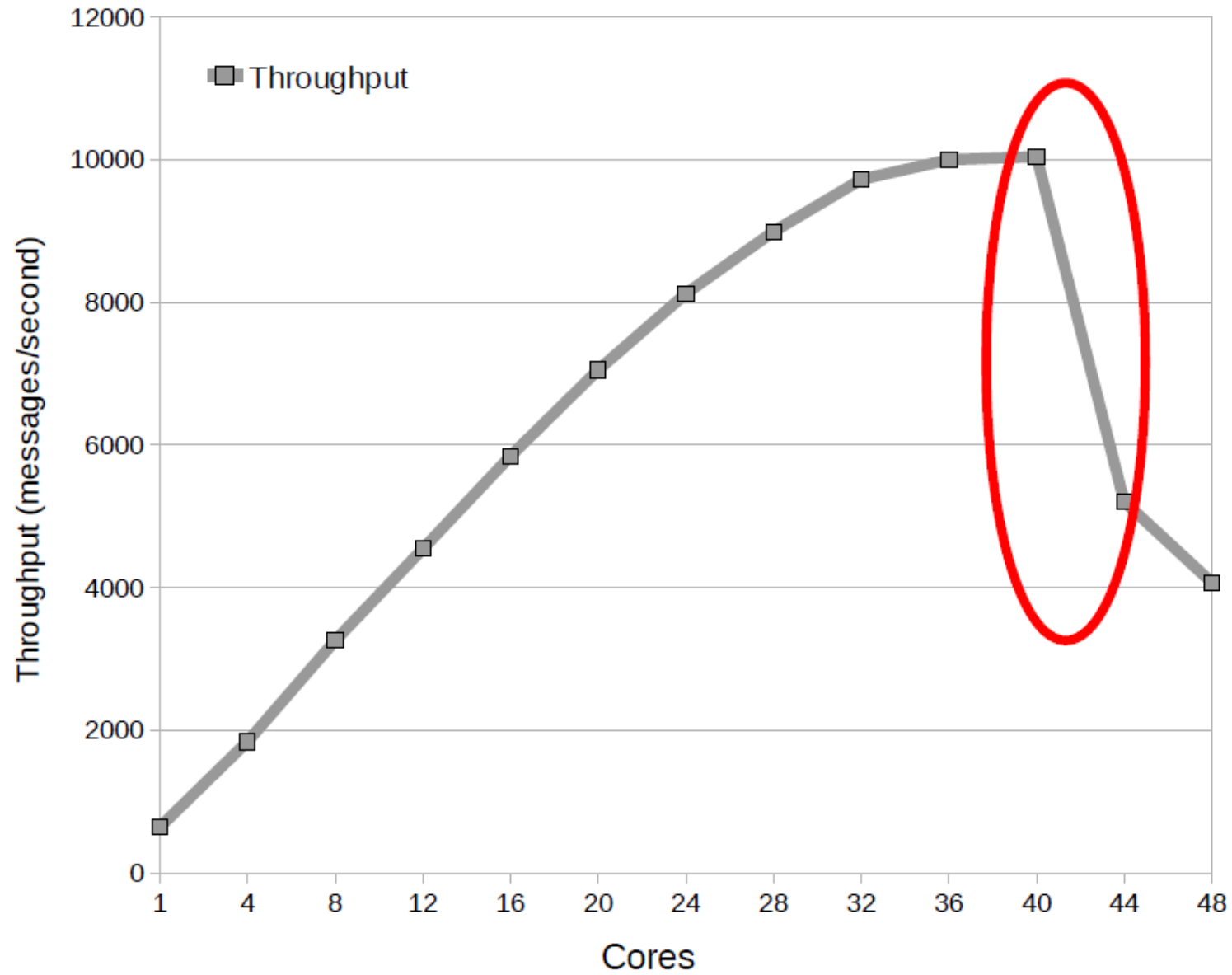
What is the main problem with locks?

- Scalability

48-core AMD server



Exim collapse



Oprofile results

40 cores:
10000 msg/sec

samples	%	app name	symbol name
2616	7.3522	vmlinux	radix_tree_lookup_slot
2329	6.5456	vmlinux	unmap_vmas
2197	6.1746	vmlinux	filemap_fault
1488	4.1820	vmlinux	__do_fault
1348	3.7885	vmlinux	copy_page_c
1182	3.3220	vmlinux	unlock_page
966	2.7149	vmlinux	page_fault

48 cores:
4000 msg/sec

samples	%	app name	symbol name
13515	34.8657	vmlinux	lookup_mnt
2002	5.1647	vmlinux	radix_tree_lookup_slot
1661	4.2850	vmlinux	filemap_fault
1497	3.8619	vmlinux	unmap_vmas
1026	2.6469	vmlinux	__do_fault
914	2.3579	vmlinux	atomic_dec
896	2.3115	vmlinux	unlock_page

Exim collapse

- `sys_open` eventually calls:

```
struct vfsmount *lookup_mnt(struct path *path)
{
    struct vfsmount *mnt;
    spin_lock(&vfsmount_lock);
    mnt = hash_get(mnts, path);
    spin_unlock(&vfsmount_lock);
    return mnt;
}
```

Exim collapse

- sys_open eventually calls:

```
struct vfsmount *lookup_mnt(struct path *path)
{
    struct vfsmount *mnt;
    spin_lock(&vfsmount_lock);
    mnt = hash_get(mnts, path);
    spin_unlock(&vfsmount_lock);
    return mnt;
}
```

← Critical section is short. Why does it cause a scalability bottleneck?

- spin_lock and spin_unlock use many more cycles than the critical section

Ticket lock in Linux

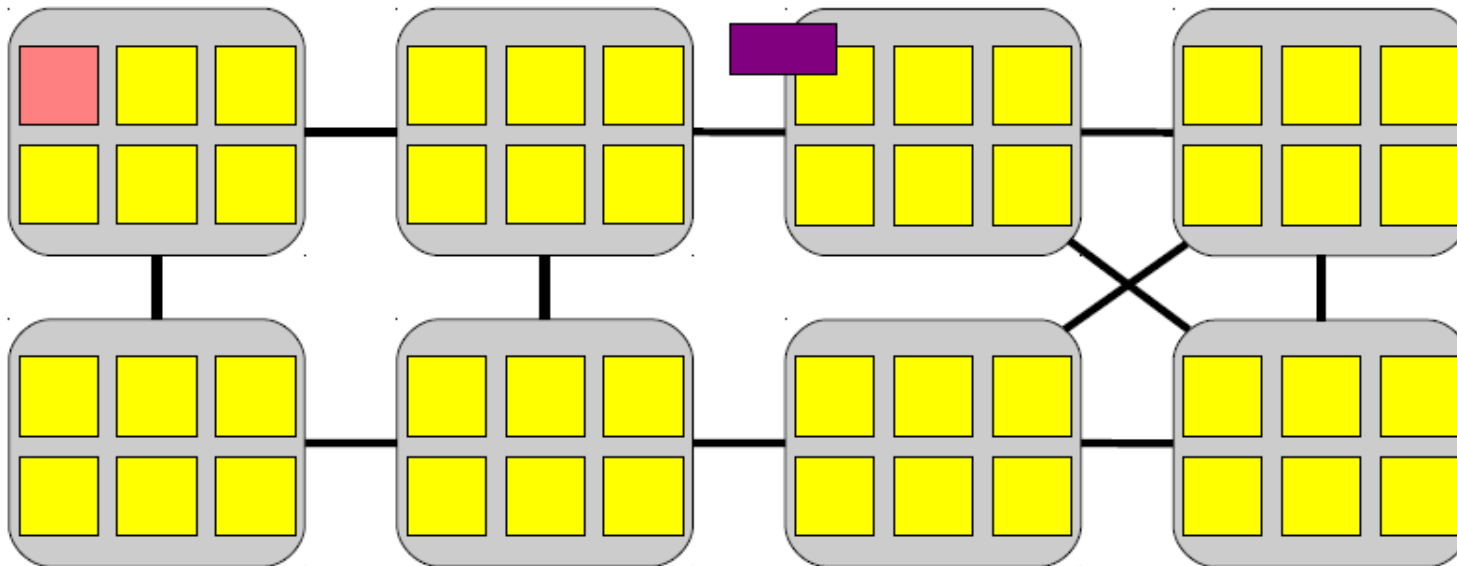
```
struct spinlock_t {  
    int current_ticket ;  
    int next_ticket ;  
}  
  
void spin_lock ( spinlock_t *lock)  
{  
    int t = atomic_fetch_and_inc (&lock -> next_ticket );  
    while (t != lock -> current_ticket )  
        ; /* spin */  
}  
  
void spin_unlock ( spinlock_t *lock)  
{  
    lock -> current_ticket ++;  
}
```

Spin lock implementation

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



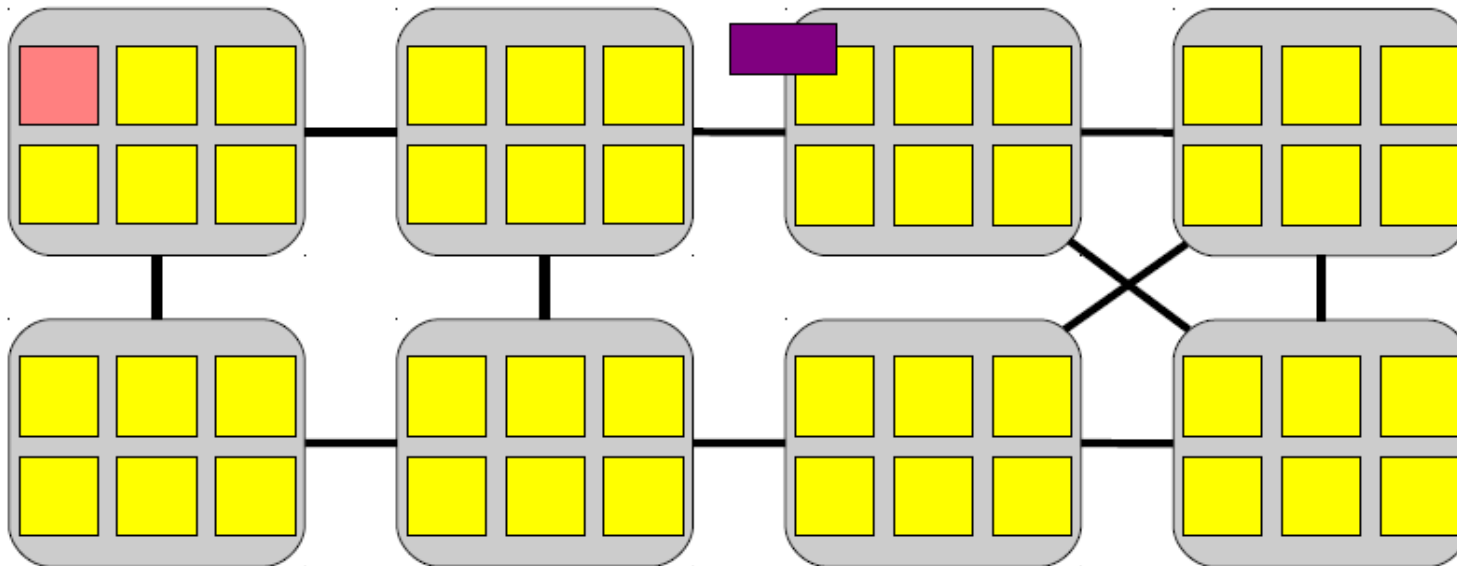
Spin lock implementation

Allocate a ticket

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



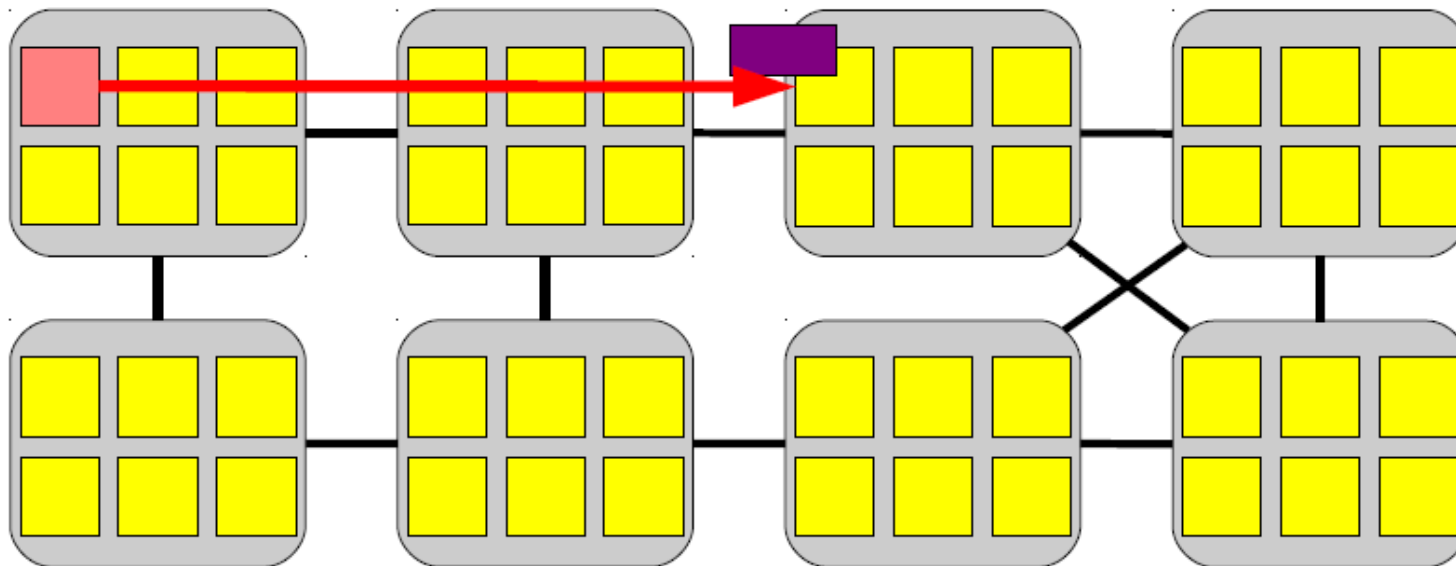
Spin lock implementation

Allocate a ticket

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



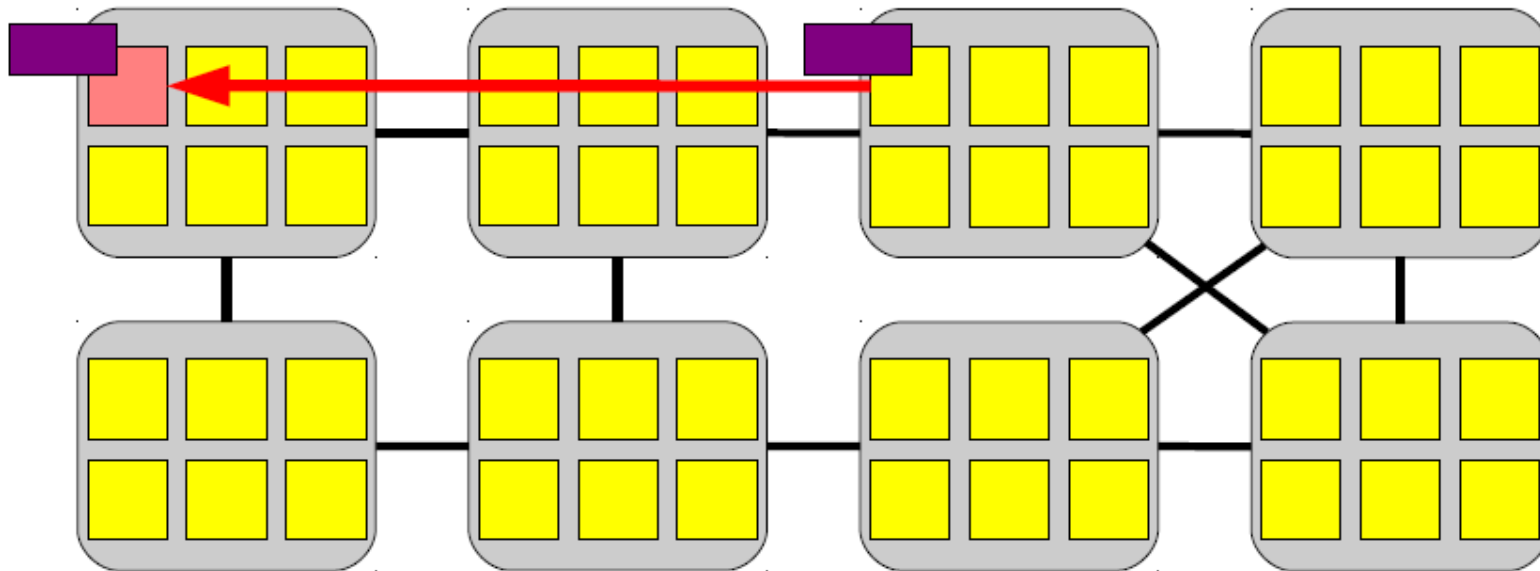
Spin lock implementation

Allocate a ticket

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



Spin lock implementation

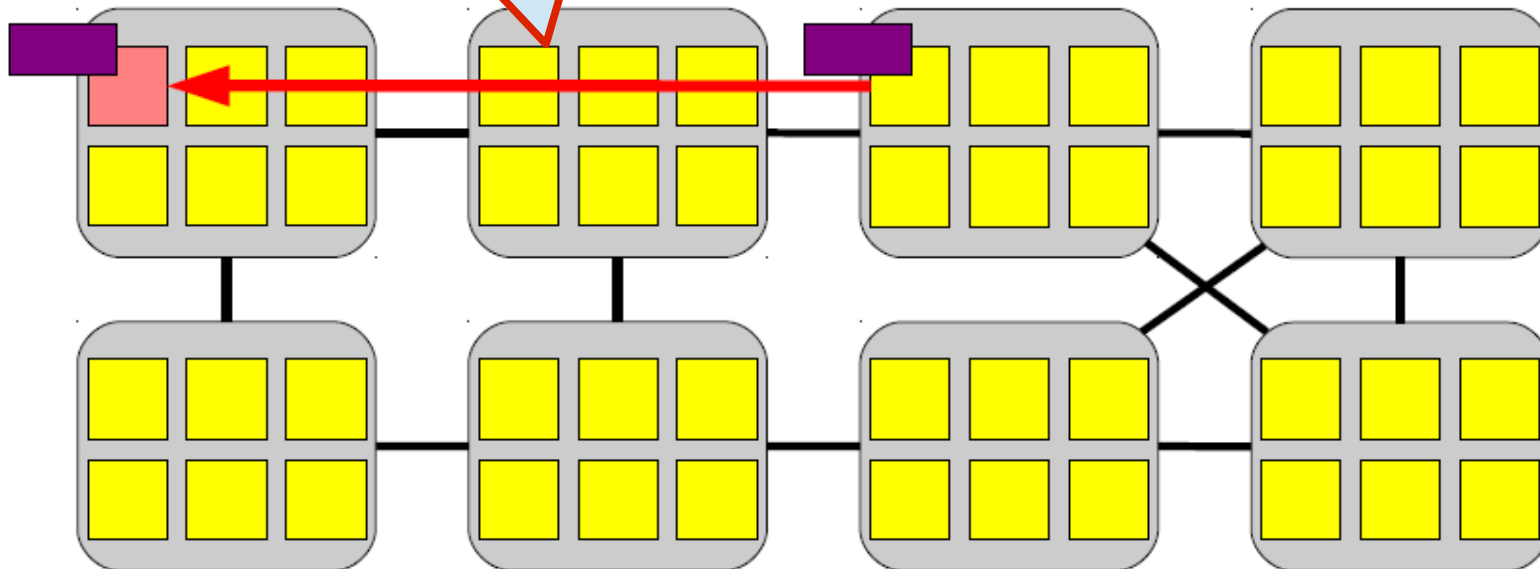
Allocate a ticket

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```

120-420 cycles



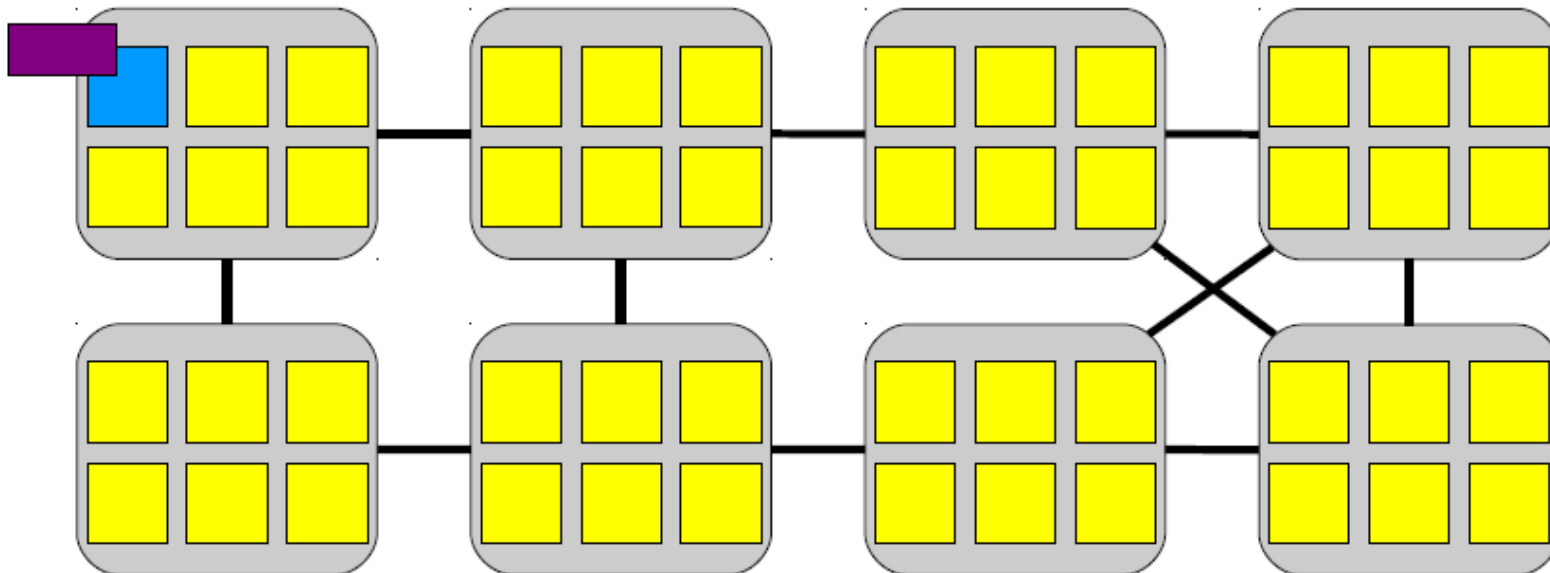
Spin lock implementation

Update the ticket value

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



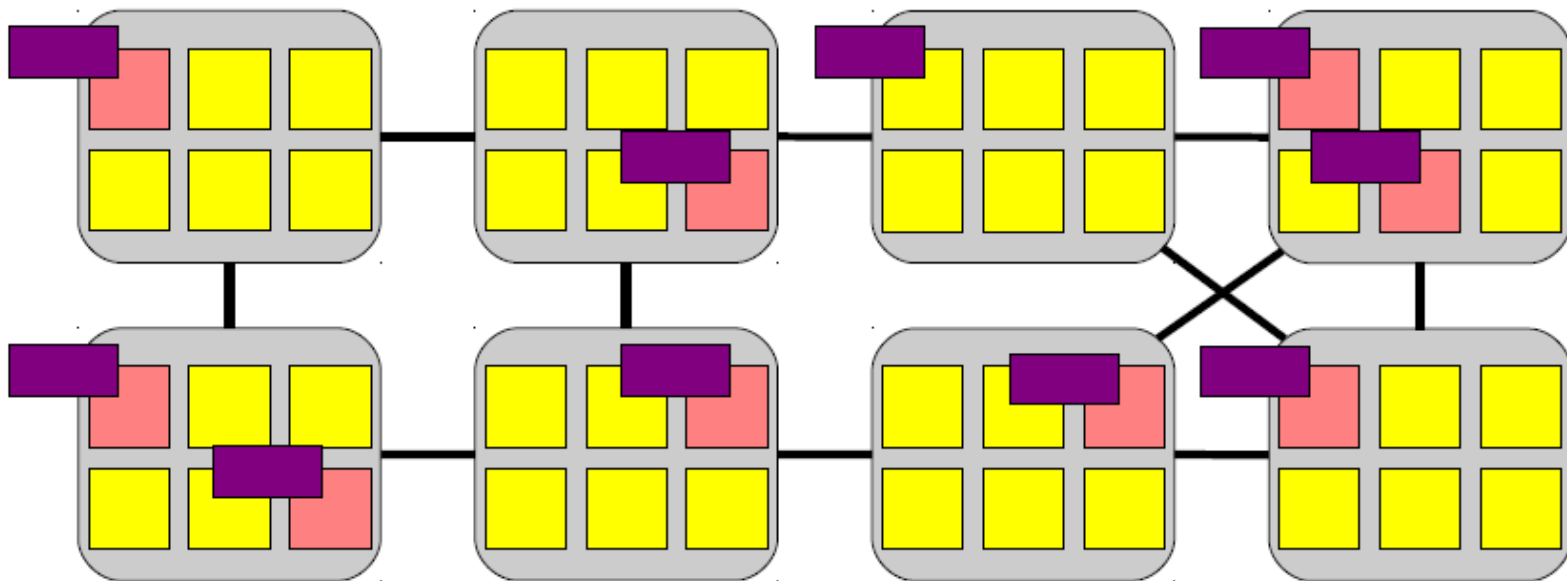
Spin lock implementation

Bunch of cores are spinning

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(&lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```



Spin lock implementation

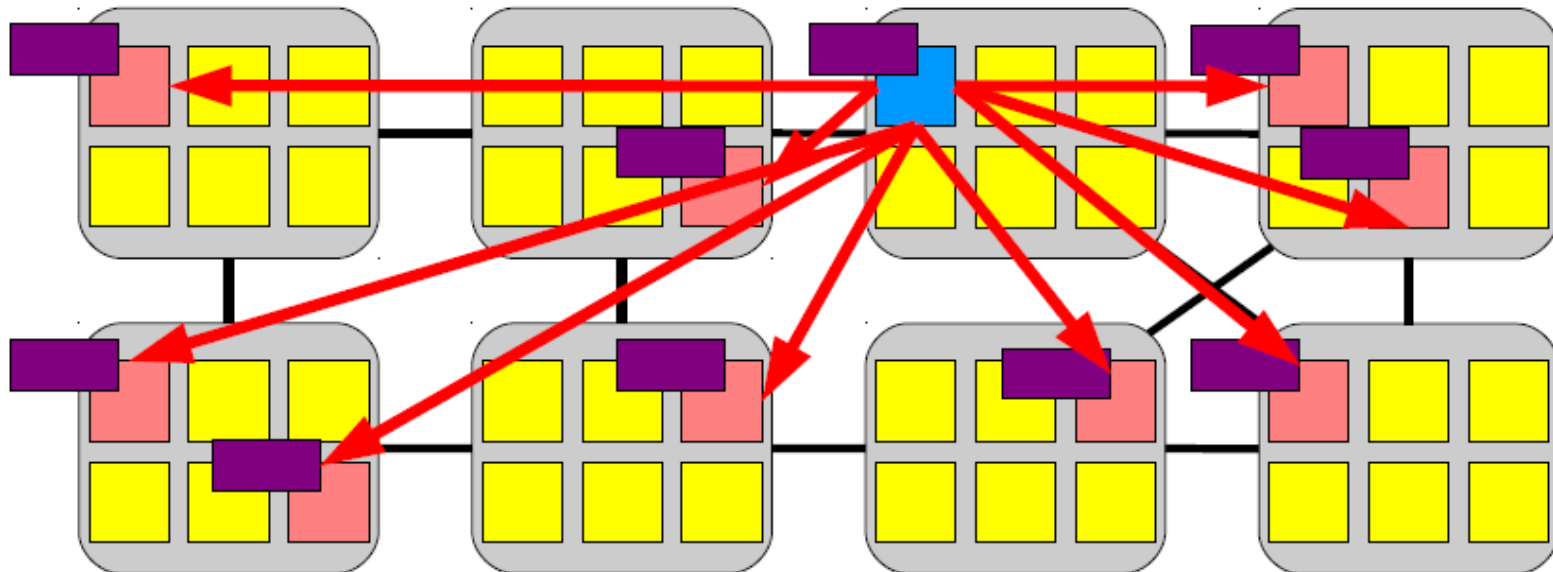
```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

Broadcast message
(invalidate the value)

```
{
    ticket;
```

```
int next_ticket;
```



Spin lock implementation

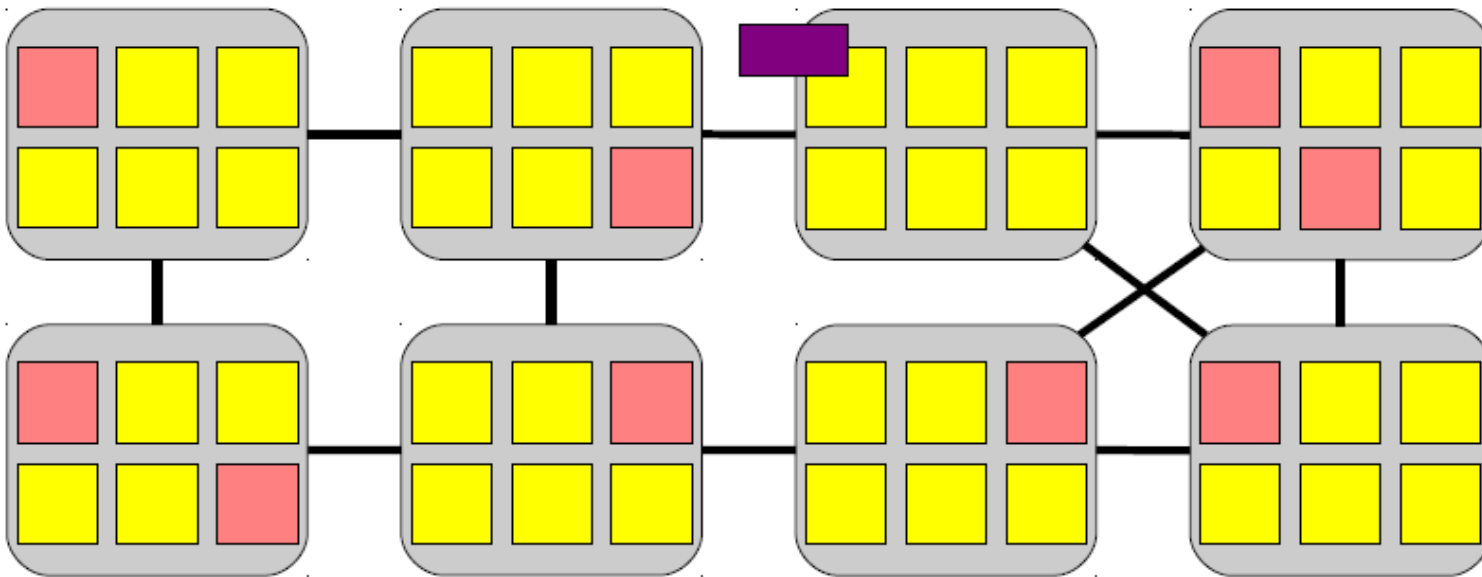
```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
```

```
    lock->current_ticket++;
```

Cores don't have the value of current_ticket

```
    lock->next_ticket++;
}
```



Spin lock implementation

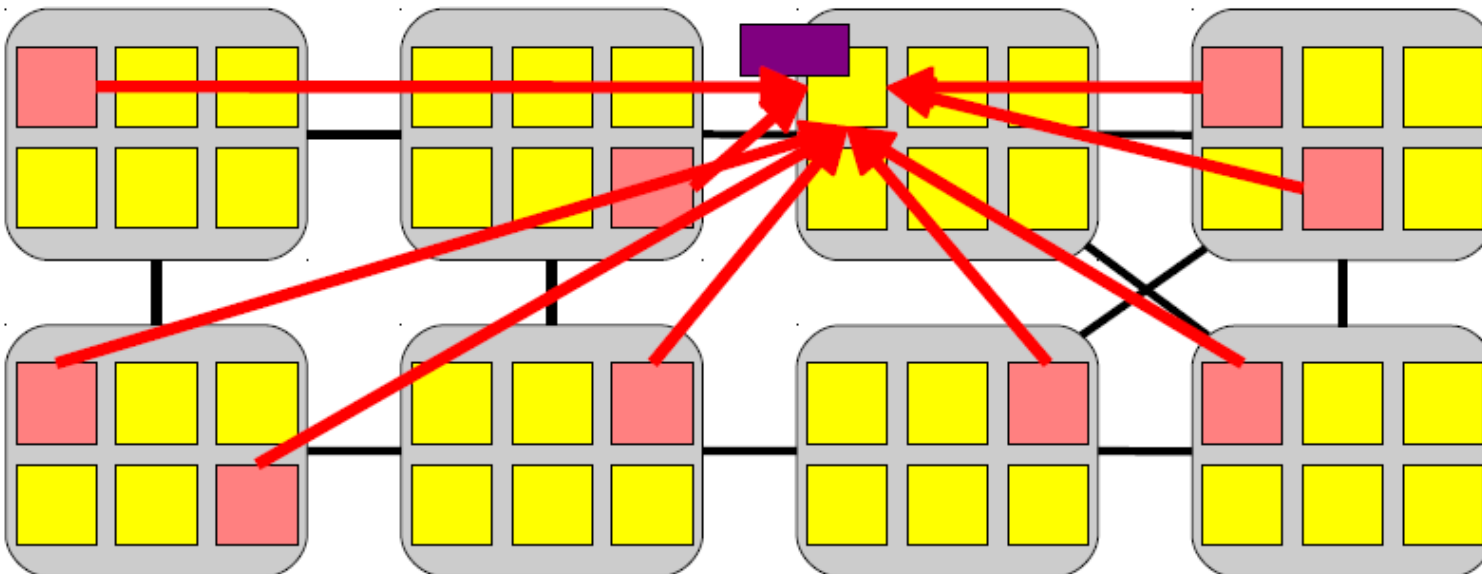
```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

Re-read the value

```
    t = atomic_inc(lock->next_ticket);
```

```
    while (t != lock->current_ticket)
        ; /* Spin */
}
```



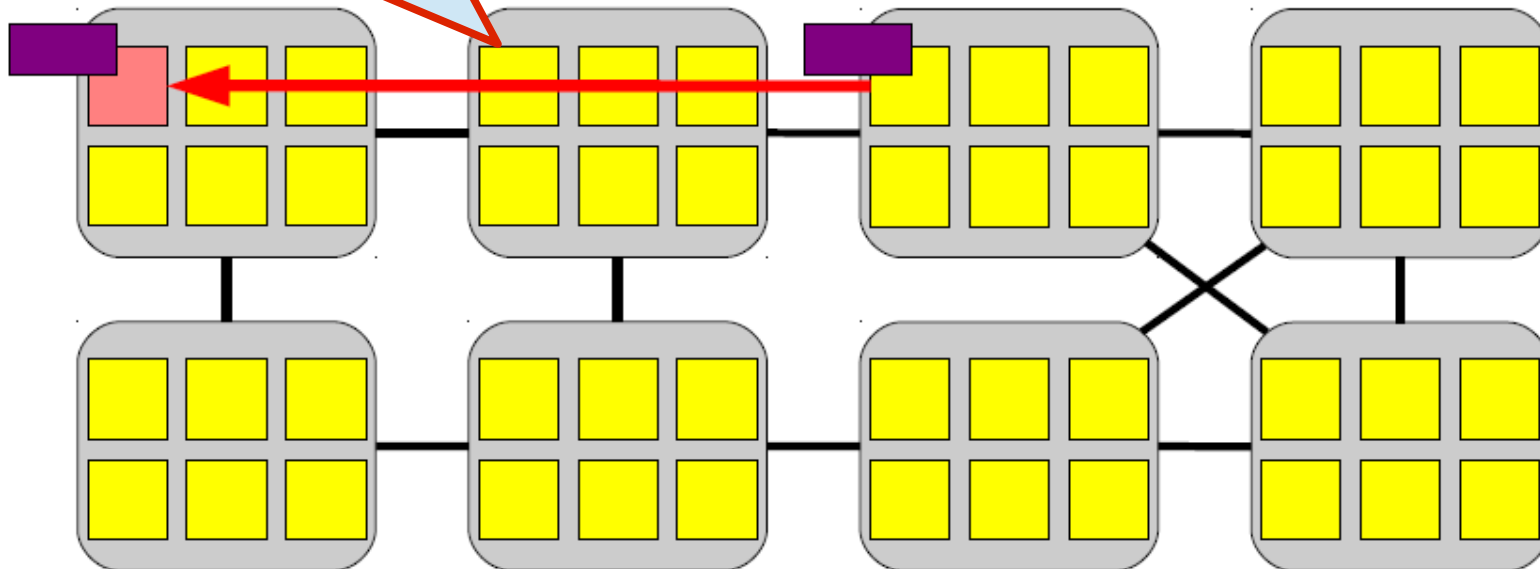
Spin lock implementation

```
void spin_lock(spinlock_t *lock)
{
    t = atomic_inc(lock->next_ticket);
    while (t != lock->current_ticket)
        ; /* Spin */
}
```

```
void spin_unlock(spinlock_t *lock)
{
    lock->current_ticket++;
}
```

```
struct spinlock_t {
    int current_ticket;
    int next_ticket;
}
```

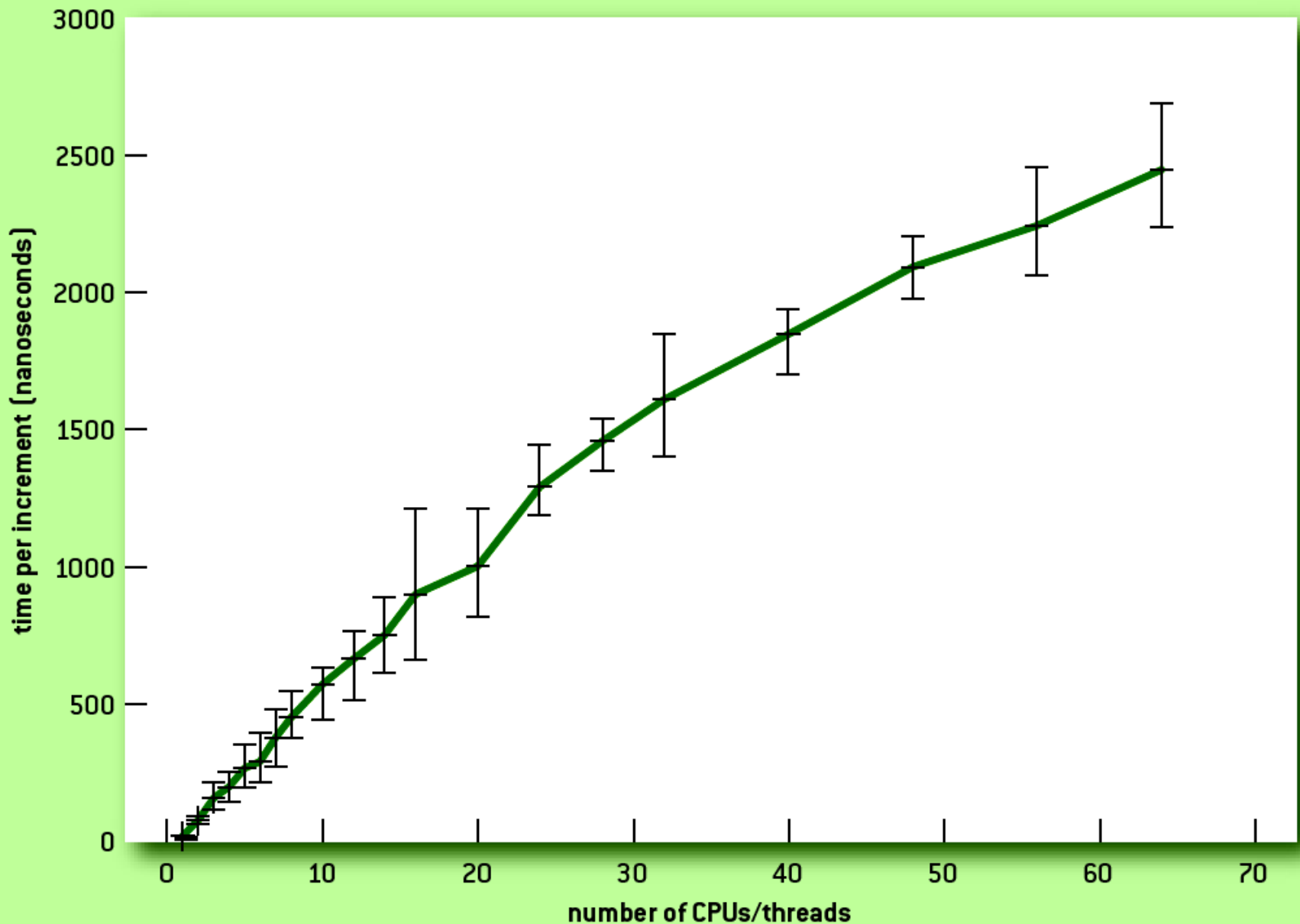
$(120-420) * N/2$ cycles



- In most architectures, the cache-coherence reads are serialized (either by a shared bus or at the cache line's home or directory node)
- Thus completing them all takes time proportional to the number of cores.
- The core that is next in line for the lock can expect to receive its copy of the cache line midway through this process.
 - $N/2$

Atomic synchronization primitives
do not scale well

Atomic increment on 64 cores



Thank you!