# Concurrent programming

# Spin Locks and Contention

Modified by Piotr Witkowski

# Focus so far: Correctness and Progress

- Models
  - Accurate (we never lied to you)
  - But idealized (so we forgot to mention a few things)
- Protocols
  - Elegant
  - Important
  - But naïve

# New Focus: Performance

- Models
  - More complicated (not the same as complex!)
  - Still focus on principles (not soon obsolete)
- Protocols
  - Elegant (in their fashion)
  - Important (why else would we pay attention)
  - And realistic (your mileage may vary)

# Kinds of Architectures

- SISD (Uniprocessor)
  - Single instruction stream
  - Single data stream
- SIMD (Vector)
  - Single instruction
  - Multiple data
- MIMD (Multiprocessors)
  - Multiple instruction
  - Multiple data.

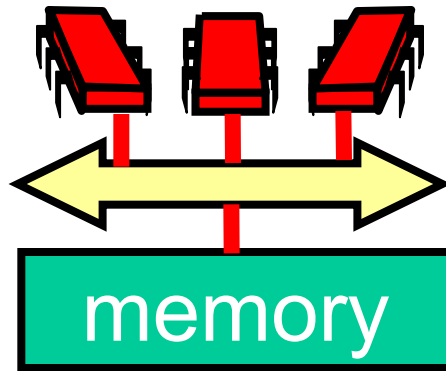# Kinds of Architectures

- SISD (Uniprocessor)
  - Single instruction stream
  - Single data stream
- SIMD (Vector)
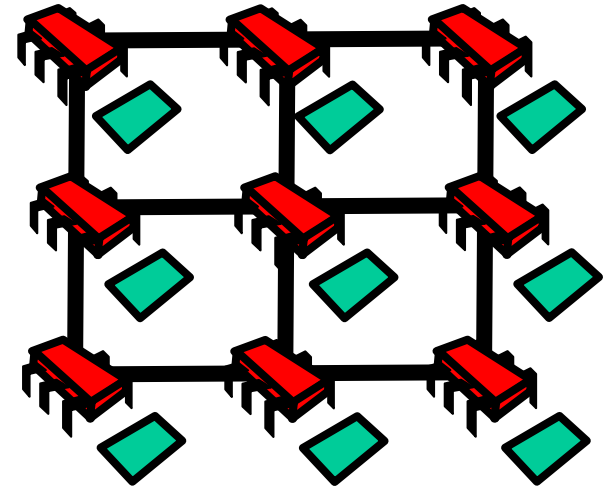  - Single instruction
  - Multiple data

Our space

- MIMD (Multiprocessors)
  - Multiple instruction
  - Multiple data.

# MIMD Architectures



**Shared Bus**

**Distributed**

- Memory Contention
- Communication Contention
- Communication Latency

# Today: Revisit Mutual Exclusion

- Performance, not just correctness
- Proper use of multiprocessor architectures
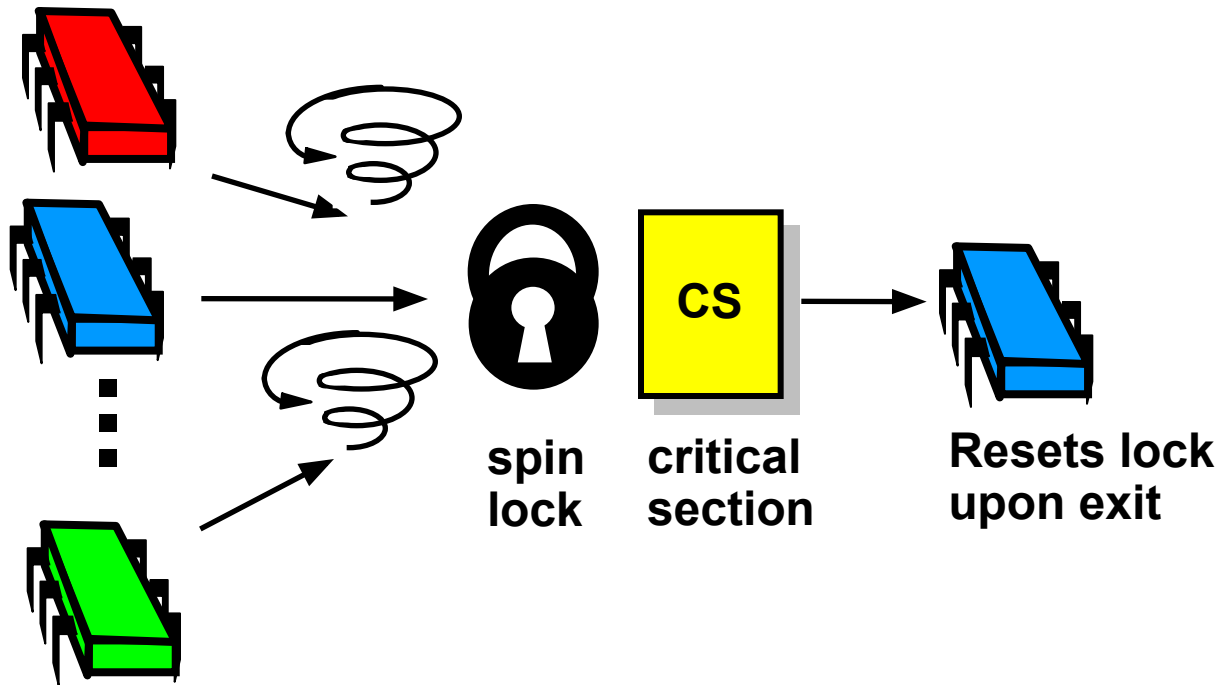- A collection of locking algorithms…

# What Should you do if you can't get a lock?

- Keep trying
  - "spin" or "busy-wait"
  - Good if delays are short
- Give up the processor
  - Good if delays are long
  - Always good on uniprocessor
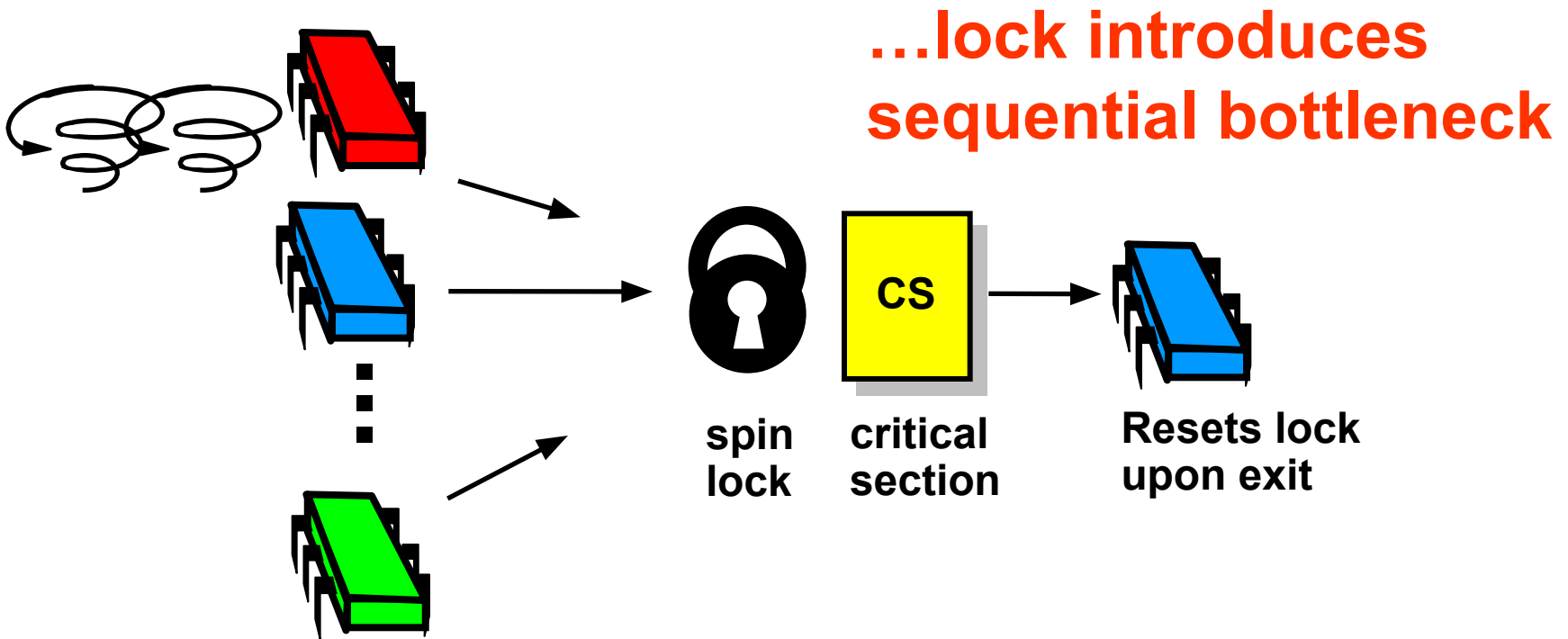
# What Should you do if you can't get a lock?

- Keep trying
  - "spin" or "busy-wait"
  - Good if delays are short
- Give up the processor
  - Good if delays are long
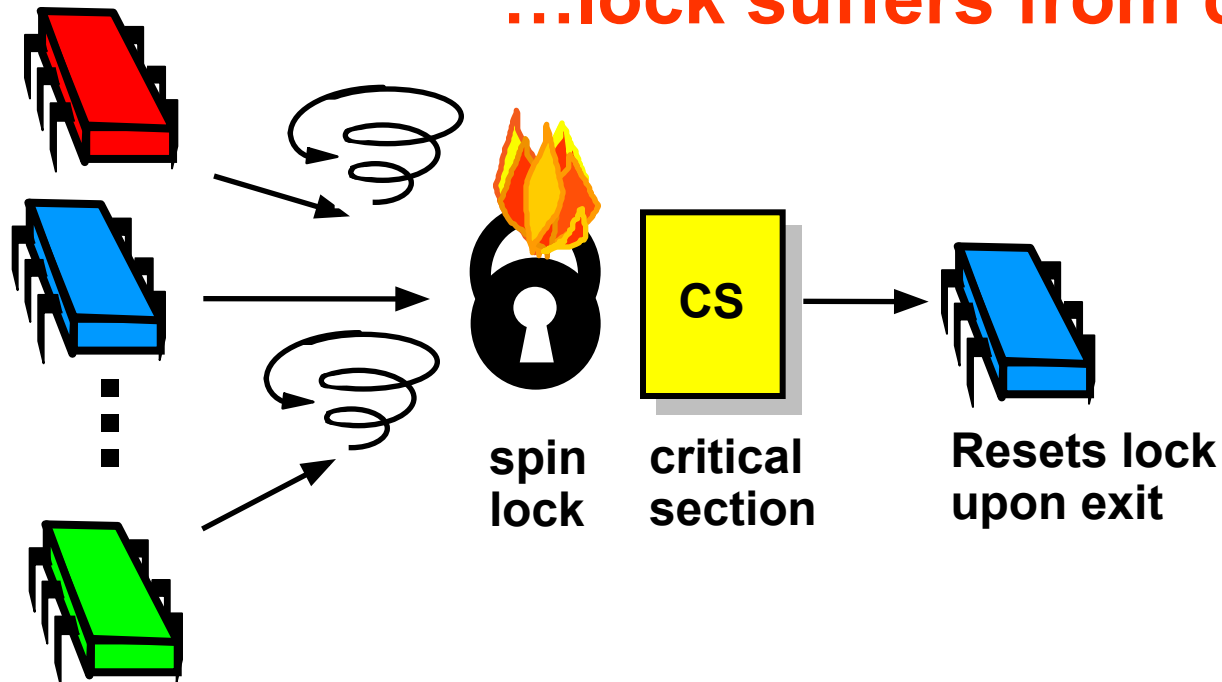  - Always good on uniprocessor

our focus

# Basic Spin-Lock



**spin lock**
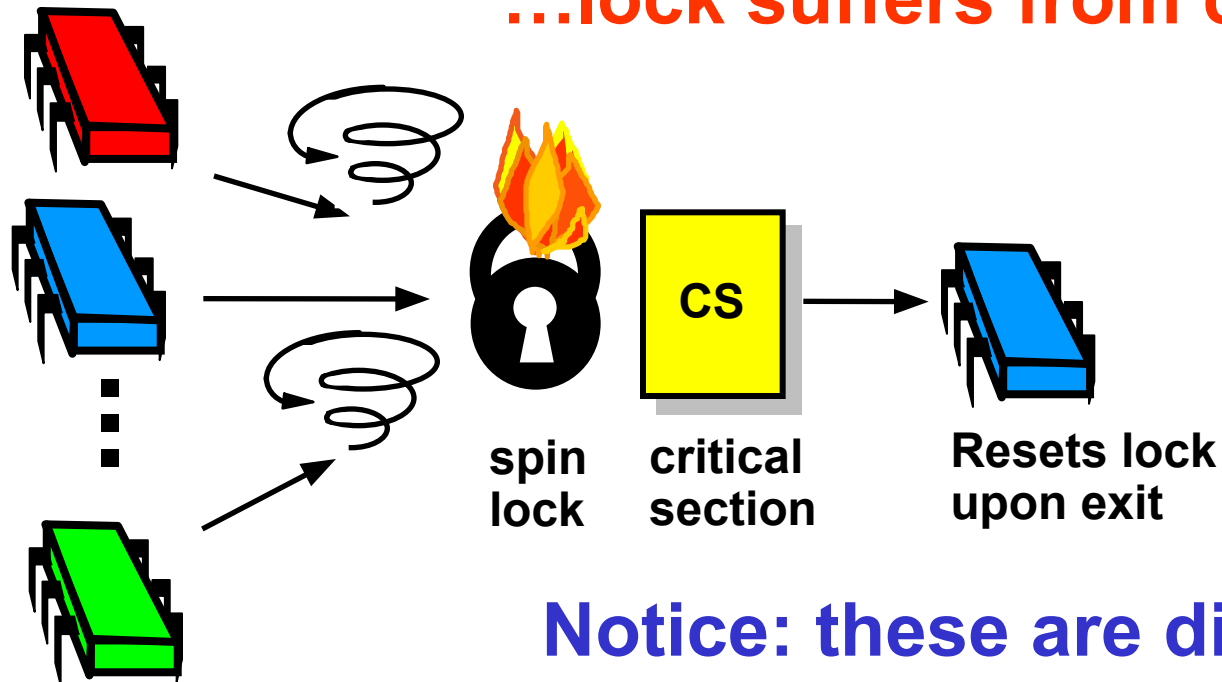
**critical section**

**Resets lock upon exit**

# Basic Spin-Lock

**…lock introduces sequential bottleneck**



**spin lock**  **critical section**  **Resets lock upon exit**

# Basic Spin-Lock

**…lock suffers from contention**



**spin lock**   **critical section**   **Resets lock upon exit**

# Basic Spin-Lock



**…lock suffers from contention**

**spin lock**

**critical section**

**Resets lock upon exit**

**Notice: these are distinct phenomena**

# Basic Spin-Lock

**…lock suffers from contention**



spin
lock

critical
section

Resets lock
upon exit

**Seq Bottleneck ⯈ no parallelism**

# Basic Spin-Lock

**…lock suffers from contention**



**spin lock**    **critical section**    **Resets lock upon exit**

**Contention 🡒 ???**

# Review: Test-and-Set

- Boolean value
- Test-and-set (TAS)
  - Swap **true** with current value
  - Return value tells if prior value was **true** or **false**
- Can reset just by writing **false**
- TAS aka "getAndSet"

# Review: Test-and-Set

```
public class AtomicBoolean {
 boolean value;

 public synchronized boolean
  getAndSet(boolean newValue) {
    boolean prior = value;
    value = newValue;
    return prior;
 }
}
```

# Review: Test-and-Set

```
public class AtomicBoolean {
  boolean value;

  public synchronized boolean
    getAndSet(boolean newValue) {
      boolean prior = value;
      value = newValue;
      return prior;
    }
}
```

**Package**
**java.util.concurrent.atomic**

# Review: Test-and-Set

```
public class AtomicBoolean {
 boolean value;

 public synchronized boolean
  getAndSet(boolean newValue) {
    boolean prior = value;
    value = newValue;
    return prior;
  }
}
```

**Swap old and new values**

# Review: Test-and-Set

```
AtomicBoolean lock
 = new AtomicBoolean(false)
…
boolean prior = lock.getAndSet(true)
```

# Review: Test-and-Set

```
AtomicBoolean lock
 = new AtomicBoolean(false)
...
boolean prior = lock.getAndSet(true)
```

**Swapping in true is called
"test-and-set" or TAS**

# Test-and-Set Locks

- Locking
  - Lock is free: value is false
  - Lock is taken: value is true
- Acquire lock by calling TAS
  - If result is false, you win
  - If result is true, you lose
- Release lock by writing false

# Test-and-set Lock

```
class TASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
  while (state.getAndSet(true)) {}
 }

 void unlock() {
  state.set(false);
}}
```

# Test-and-set Lock

```
class TASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
  while (state.getAndSet(true)) {}
 }

 void unlock() {
  state
}}
```

**Lock state is AtomicBoolean**
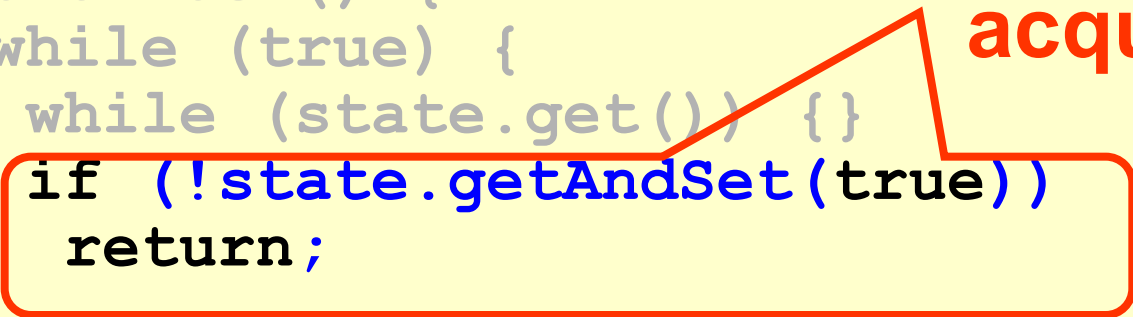
# Test-and-set Lock

```
class TASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
   while (state.getAndSet(true)) {}
 }

 void unlock() {
  sta
 }}
```

**Keep trying until lock acquired**

# Test-and-set Lock

```
class TA
 AtomicB
  new At

 void lock() {
  while (state.getAndSet(true)) {}
 }


 void unlock() {
  state.set(false);
 }}
```

**Release lock by resetting state to false**

# Space Complexity

- TAS spin-lock has small "footprint"
- N thread spin-lock uses O(1) space
- As opposed to O(n) Peterson/Bakery
- How did we overcome the $\Omega(n)$ lower bound?
- We used a RMW operation…

# Performance

- Experiment
  - *n* threads
  - Increment shared counter 1 million times
- How long should it take?
- How long does it take?

# Graph

# Mystery #1



time

threads

TAS lock

Ideal

What is going on?

# Test-and-Test-and-Set Locks

- Lurking stage
  - Wait until lock "looks" free
  - Spin while read returns true (lock taken)
- Pouncing state
  - As soon as lock "looks" available
  - Read returns false (lock free)
  - Call TAS to acquire lock
  - If TAS loses, back to lurking

# Test-and-test-and-set Lock

```
class TTASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
  while (true) {
   while (state.get()) {}
   if (!state.getAndSet(true))
     return;
  }
}
```

# Test-and-test-and-set Lock

```
class TTASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
  while (true) {
   while (state.get()) {}
   if (!state.getAndSet(true))
    return;
  }
}
```

**Wait until lock looks free**

# Test-and-test-and-set Lock

```
class TTASlock {
 AtomicBoolean state =
  new AtomicBoolean(false);

 void lock() {
  while (true) {
   while (state.get()) {}
    if (!state.getAndSet(true))
     return;
  }
 }
}
```

**Then try to acquire it**

# Mystery #2



time

threads

TAS lock

TTAS lock

Ideal

# Mystery

- Both
  - TAS and TTAS
  - Do the same thing (in our model)
- Except that
  - TTAS performs much better than TAS
  - Neither approaches ideal

# Opinion

- Our memory abstraction is broken
- TAS & TTAS methods
  - Are provably the same (in our model)
  - Except they aren't (in field tests)
- Need a more detailed model …

# Bus-Based Architectures



cache    cache    cache

Bus

memory

# Bus-Based Architectures



Random access memory
(10s of cycles)

cache

memory

# Bus-Based Architectures

Shared Bus
- Broadcast medium
- One broadcaster at a time
- Processors and memory all "snoop"

cache     cache     cache

Bus

memory

Bus-

Per-Processor Caches
• Small
• Fast: 1 or 2 cycles
• Address & state information

cache          cache          cache

Bus

memory

# Granularity

- Caches operate at a larger granularity than a word

- Cache line: fixed-size block containing the address (today 64 or 128 bytes)

# Locality

- If you use an address now, you will probably use it again soon
  - Fetch from cache, not memory
- If you use an address now, you will probably use a nearby address soon
  - In the same cache line

# L1 and L2 Caches



**L2**

**L1**

# L1 and L2 Caches

**L2**

**L1**

**Small & fast
1 or 2 cycles**

# L1 and L2 Caches

**Larger and slower**
**10s of cycles**
**~128 byte line**

**L2**

**L1**

# Jargon Watch

- Cache hit
  - "I found what I wanted in my cache"
  - Good Thing™

# Jargon Watch

- Cache hit
  - "I found what I wanted in my cache"
  - Good Thing™
- Cache miss
  - "I had to shlep all the way to memory for that data"
  - Bad Thing™

# Cave Canem

- This model is still a simplification
  - But not in any essential way
  - Illustrates basic principles
- Will discuss complexities later

# Fully Associative Cache

- Any line can be anywhere in the cache
  - Advantage: can replace any line
  - Disadvantage: hard to find lines

# Direct Mapped Cache

- Every address has exactly 1 slot
  - Advantage: easy to find a line
  - Disadvantage: must replace fixed line

# K-way Set Associative Cache

- Each slot holds k lines
  - Advantage: pretty easy to find a line
  - Advantage: some choice in replacing line

# Multicore Set Associativity

- k is 8 or even 16 and growing…
  - Why? Because cores share sets
  - Threads cut effective size if accessing different data

# Cache Coherence

- A and B both cache address x

- A writes to x

  – Updates cache

- How does B find out?

- Many cache coherence protocols in literature

# MESI

- Modified
  - Have modified cached data, must write back to memory

# MESI

- Modified
  - Have modified cached data, must write back to memory
- Exclusive
  - Not modified, I have only copy

# MESI

- Modified
  - Have modified cached data, must write back to memory
- Exclusive
  - Not modified, I have only copy
- Shared
  - Not modified, may be cached elsewhere

# MESI

- Modified
  - Have modified cached data, must write back to memory
- Exclusive
  - Not modified, I have only copy
- Shared
  - Not modified, may be cached elsewhere
- Invalid
  - Cache contents not meaningful

# Processor Issues Load Request

load x

cache cache cache

Bus

memory data

# Memory Responds



E

cache    cache    cache

Bus

Got it!

memory    data

# Processor Issues Load Request



Load x

E

data

cache

cache

Bus

memory    data

# Other Processor Responds

# Modify Cached Data

# Write-Through Cache

Write x!

S data   S d a   cache

Bus

memory data

# Write-Through Caches

- Immediately broadcast changes
- Good
  – Memory, caches always agree
  – More read hits, maybe
- Bad
  – Bus traffic on all writes
  – Most writes to unshared data
  – For example, loop indexes …

# Write-Through Caches

- Immediately broadcast changes
- Good
  - Memory, caches always agree
  - More read hits, maybe
- Bad
  - Bus traffic on all writes
  - Most writes to unshared data
  - For example, loop indexes …

"show stoppers"

# Write-Back Caches

- Accumulate changes in cache
- Write back when line evicted
  - Need the cache for something else
  - Another processor wants it

# Invalidate



I

M

# Invalidate



This cache acquires write permission

# Invalidate

Other caches lose read permission



cache          data          cache

This cache acquires write permission

# Invalidate

Memory provides data only if not present in any cache, so no need to change it now (expensive)

Bus

memory data

# Mutual Exclusion

- What do we want to optimize?
  - Bus bandwidth used by spinning threads
  - Release/Acquire latency
  - Acquire latency for idle lock

# Simple TASLock

- TAS invalidates cache lines

- Spinners
  - Miss in cache
  - Go to bus

- Thread wants to release lock
  - delayed behind spinners

# Test-and-test-and-set

- Wait until lock "looks" free
  - Spin on local cache
  - No bus use while lock busy
- Problem: when lock is released
  - Invalidation storm …

# Local Spinning while Lock is Busy

# On Release

invalid

invalid

free

Bus

memory free

# On Release

Everyone misses, rereads



miss

miss

free

Bus

memory  free

# On Release

Everyone tries TAS



TAS(…)   TAS(…)   free

Bus

memory   free

# Problems

- Everyone misses
  - Reads satisfied sequentially
- Everyone does TAS
  - Invalidates others' caches
- Eventually quiesces after lock acquired
  - How long does this take?

# Measuring Quiescence Time



- Acquire lock
- Pause without using bus
- Use bus heavily

If pause > quiescence time,
critical section duration independent of number of threads

If pause < quiescence time,
critical section duration slower with more threads

# Quiescence Time



**Increses linearly with the number of processors for bus architecture**

time

threads

# Mystery Explained

**TAS lock**

**TTAS lock**

**Ideal**

tim

thread
s

Better than TAS but still not as good as ideal

# Solution: Introduce Delay

- **If the lock looks free**
  - **But I fail to get it**
- **There must be contention**
  - **Better to back off than to collide again**

time ----- |---------|-------|-------|------|

$r_2d$    $r_1d$    $d$

spin lock

# Dynamic Example: Exponential Backoff

time

**4d**        **2d**        **d**

spin lock

If I fail to get lock
- Wait random duration before retry
- Each subsequent failure doubles expected wait

# Exponential Backoff Lock

```java
public class Backoff implements lock {
 public void lock() {
   int delay = MIN_DELAY;
   while (true) {
    while (state.get()) {}
    if (!lock.getAndSet(true))
      return;
    sleep(random() % delay);
    if (delay < MAX_DELAY)
     delay = 2 * delay;
 }}}
```

# Exponential Backoff Lock

```
public class Backoff implements lock {
 public void lock() {
   int delay = MIN_DELAY;
   while (true) {
    while (state.get()) {}
    if (!lock.getAndSet(true))
     return;
    sleep(random() % delay);
    if (delay < MAX_DELAY)
     delay = 2 * delay;
 }}}
```

**Fix minimum delay**

# Exponential Backoff Lock

```
public class Backoff implements lock {
 public void lock() {
   int delay = MIN_DELAY;
   while (true) {
    while (state.get()) {}
    if (!lock.getAndSet(true))
     return;
    sleep(random() % delay);
    if (delay < MAX_DELAY)
     delay = 2 * delay;
}}}
```

**Wait until lock looks free**

# Exponential Backoff Lock

```
public class Backoff implements lock {
 public void lock() {
   int delay = MIN_DELAY;
   while (true) {
    while (state.get()) {}
    if (!lock.getAndSet(true))
     return;
    sleep(random() % delay);
    if (delay < MAX_DELAY)
     delay = 2 * delay;
 }}}
```

**If we win, return**

# Exponential Backoff Lock

```
public class Backoff implements lock {
 public
  int delay = MIN_DELAY;
  while (true) {
   while (state.get()) {}
   if (!lock.getAndSet(true))
    return;
   sleep(random() % delay);
   if (delay < MAX_DELAY)
    delay = 2 * delay;
 }}}
```

**Back off for random duration**

# Exponential Backoff Lock

```
public class Backoff implements lock {
  public
  int delay = MIN_DELAY;
  while (true) {
    while (state.get()) {}
    if (!lock.getAndSet(true))
      return;
    sleep(random() % delay);
    if (delay < MAX_DELAY)
      delay = 2 * delay;
}}}
```

**Double max delay, within reason**

# Spin-Waiting Overhead



time

threads

TTAS Lock

Backoff lock

# Backoff: Other Issues

- Good
  - Easy to implement
  - Beats TTAS lock
- Bad
  - Must choose parameters carefully
  - Not portable across platforms

# Actual Data on 40-Core Machine



Lock Scalability - Latency

# Idea

- Avoid useless invalidations
  - By keeping a queue of threads
- Each thread
  - Notifies next in line
  - Without bothering the others

# Anderson Queue Lock

idle

next

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

acquiring

next

getAndIncrement

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

acquiring

next

getAndIncrement

flags

| T | F | F | F | F | F | F | F |

# Anderson Queue Lock

acquired

next

Mine!

flags

| T | F | F | F | F | F | F | F |

# Anderson Queue Lock

acquired        acquiring

next

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

next
acquired   acquiring

getAndIncrement

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

acquired        acquiring

next



getAndIncrement

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

acquired     acquiring

next

flags

| T | F | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

released     acquired

next

flags

| T | T | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

# Anderson Queue Lock

next

released     acquired

flags

| T | T | F | F | F | F | F | F |
|---|---|---|---|---|---|---|---|

Yow!

# Anderson Queue Lock

```
class ALock implements Lock {
 boolean[] flags={true,false,…,false};
 AtomicInteger next
  = new AtomicInteger(0);
 ThreadLocal<Integer> mySlot;
```

# Anderson Queue Lock

```
class ALock implements Lock {
boolean[] flags={true,false,…,false};
AtomicInteger next
  = new AtomicInteger(0);
ThreadLocal<Integer> mySlot;
```

**One flag per thread**

# Anderson Queue Lock

```
class ALock implements Lock {
  boolean[] flags={true,false,…,false};
  AtomicInteger next
    = new AtomicInteger(0);
  ThreadLocal<Integer> mySlot;
```

**Next flag to use**

# Anderson Queue Lock

```
class ALock implements Lock {
 boolean[] flags={true,false,…,false};
 AtomicInteger next
  = new AtomicInteger(0);
 ThreadLocal<Integer> mySlot;
```

**Thread-local variable**

# Anderson Queue Lock

```
public lock() {
 mySlot = next.getAndIncrement();
 while (!flags[mySlot % n]) {};
 flags[mySlot % n] = false;
}


public unlock() {
 flags[(mySlot+1) % n] = true;
}
```

# Anderson Queue Lock

```
public lock() {
  mySlot = next.getAndIncrement();
  while (!flags[mySlot % n]) {};
  flags[mySlot % n] = false;
}

public unlock() {
  flags[(mySlot+1) % n]
}
```

**Take next slot**

# Anderson Queue Lock

```
public lock() {
 mySlot = next.getAndIncrement();
 while (!flags[mySlot % n]) {};
 flags[mySlot % n] = false;
}

public unlock() {
 flags[(mySlot+1` ^ ` ` `
}
```

**Spin until told to go**

# Anderson Queue Lock

```
public lock() {
 myslot = next.getAndIncrement();
 while (!flags[myslot % n]) {};
 flags[myslot % n] = false;
}

public unlock() {
 flags[(myslot+1) % n] = true;
}
```

**Prepare slot for re-use**

# Anderson Queue Lock

```
public lock() {
 mySlot = next.getAndIncrement();
 while (!flags[mySlot % n]) {};
 flags[mySlot % n] = false;
}


public unlock() {
 flags[(mySlot+1) % n] = true;
}
```

**Tell next thread to go**

# Local Spinning

# False Sharing



released

acquired

next

Spin on my

Result: contention

T | F | F | F | F

Spinning thread gets cache invalidation on account of store by threads it is not waiting for

Line 1

Line 2

# The Solution: Padding

# Performance



TTAS

queue

- **Shorter handover than backoff**
- **Curve is practically flat**
- **Scalable performance**

# Anderson Queue Lock

Good

- – First truly scalable lock
- – Simple, easy to implement
- – Back to FCFS order (like Bakery)

# Anderson Queue Lock

Bad

– Space hog…

– One bit per thread ◻ one cache line per thread

- What if unknown number of threads?
- What if small number of actual contenders?

# CLH Lock

- FCFS order
- Small, constant-size overhead per thread

# Initially

idle



tail



false

# Initially

idle

tail

false

Queue tail

# Initially

idle

tail

false

Lock is free

# Initially

idle

tail

false

# Purple Wants the Lock

acquiring

tail

false

# Purple Wants the Lock

acquiring



tail

false

true

# Purple Wants the Lock

acquiring



Swap

tail

false

true

# Purple Has the Lock

acquired

tail

false

true

# Red Wants the Lock

acquired

acquiring



tail

false

true

true

# Red Wants the Lock

acquired      acquiring



Swap

tail

false     true     true

# Red Wants the Lock

acquired                    acquiring



tail

false          true          true

# Red Wants the Lock

acquired          acquiring



tail

false          true          true

# Red Wants the Lock

acquired

acquiring

**Implicit Linked list**

tail

false

true

true

# Red Wants the Lock

acquired               acquiring

tail

| false | true | true |

# Red Wants the Lock

# Purple Releases

# Purple Releases

released        acquired



tail

true

# Space Usage

- Let
  - L = number of locks
  - N = number of threads
- ALock
  - $O(LN)$
- CLH lock
  - $O(L+N)$

# CLH Queue Lock

```
class QNode {
 AtomicBoolean locked =
    new AtomicBoolean(true);
}
```

# CLH Queue Lock

```
class QNode {
  AtomicBoolean locked =
    new AtomicBoolean(true);
}
```

**Not released yet**

# CLH Queue Lock

```
class CLHLock implements Lock {
 AtomicReference<QNode> tail;
 ThreadLocal<QNode> myNode
    = new QNode();
 public void lock() {
  QNode pred
    = tail.getAndSet(myNode);
  while (pred.locked) {}
}}
```

# CLH Queue Lock

```
class CLHLock implements Lock {
 AtomicReference<QNode> tail;
 ThreadLocal<QNode> myNode
    = new QNode();
 public void lock() {
  QNode pred
    = tail.getAndSet(myNode);
  while (pred.locked) {}
}}
```

**Queue tail**

# CLH Queue Lock

```
class CLHLock implements Lock {
 AtomicReference<QNode> tail;
 ThreadLocal<QNode> myNode
    = new QNode();
 public void lock() {
  QNode pred
    = tail.getAndSet(myNode);
  while (pred.locked) {}
}}
```

Thread-local QNode

# CLH Queue Lock

```
class CLHLock implements Lock {
 AtomicReference<QNode> tail;
 ThreadLocal<QNode> myNode
    = new QNode();
 public void lock() {
  QNode pred
    = tail.getAndSet(myNode);
  while (pred.locked) {}
}}
```

Swap in my node

# CLH Queue Lock

```
class CLHLock implements Lock {
 AtomicReference<QNode> tail;
 ThreadLocal<QNode> myNode
    = new QNode();
 public void lock() {
  QNode pred
    = tail.getAndSet(myNode);
  while (pred.locked) {}
}}
```

Spin until predecessor releases lock

# CLH Queue Lock

```
Class CLHLock implements Lock {
 …
 public void unlock() {
  myNode.locked.set(false);
  myNode = pred;
 }
}
```

# CLH Queue Lock

```
Class CLHLock implements Lock {

  …

  public void unlock() {
    myNode.locked.set(false);
    myNode = pred;
  }

}
```

Notify successor

# CLH Queue Lock

```
Class CLHLock implements Lock {

  …

  public void unlock() {
    myNode.locked.set(false);
    myNode = pred;
  }

}
```

Recycle predecessor's node

# CLH Queue Lock

```
Class CLHLock implements Lock {
 …
 public void unlock() {
  myNode.locked.set(false);
  myNode = pred;
 }
}
```

**myNode = pred;**

**(we don't actually reuse myNode.
Code in book shows how it's done.)**

# CLH Lock

- Good
  - Lock release affects predecessor only
  - Small, constant-sized space
- Bad
  - Doesn't work for uncached NUMA architectures

# NUMA and cc-NUMA Architectures

- Acronym:
  - **N**on-**U**niform **M**emory **A**rchitecture
  - ccNUMA = cache coherent NUMA
- Illusion:
  - Flat shared memory
- Truth:
  - No caches (sometimes)
  - Some memory regions faster than others

# NUMA Machines



**Spinning on local memory is fast**

# NUMA Machines



**Spinning on remote memory is slow**

# CLH Lock

- Each thread spins on predecessor's memory
- Could be far away ...

# MCS Lock

- FCFS order
- Spin on local memory only
- Small, Constant-size overhead

# Initially

idle



tail

false

# Acquiring

acquiring

**(allocate QNode)**

true →

tail

false →

# Acquiring

acquired

swap

true

tail

false

# Acquiring

acquired



true

tail

false

# Acquired

acquired

false

true

tail

# Acquiring

acquired

acquiring



false

tail

swap

true

# Acquiring

acquired

acquiring

false

true

tail

# Acquiring

acquired

acquiring

false

tail

true

# Acquiring

acquiring

acquired

false

true

tail

# Acquiring

acquiring

acquired

tail

true

false

# Acquiring

acquired

acquiring

true

tail

false

Yes!

# MCS Queue Lock

```
class QNode {
 volatile boolean locked = false;
 volatile qnode    next    = null;
}
```
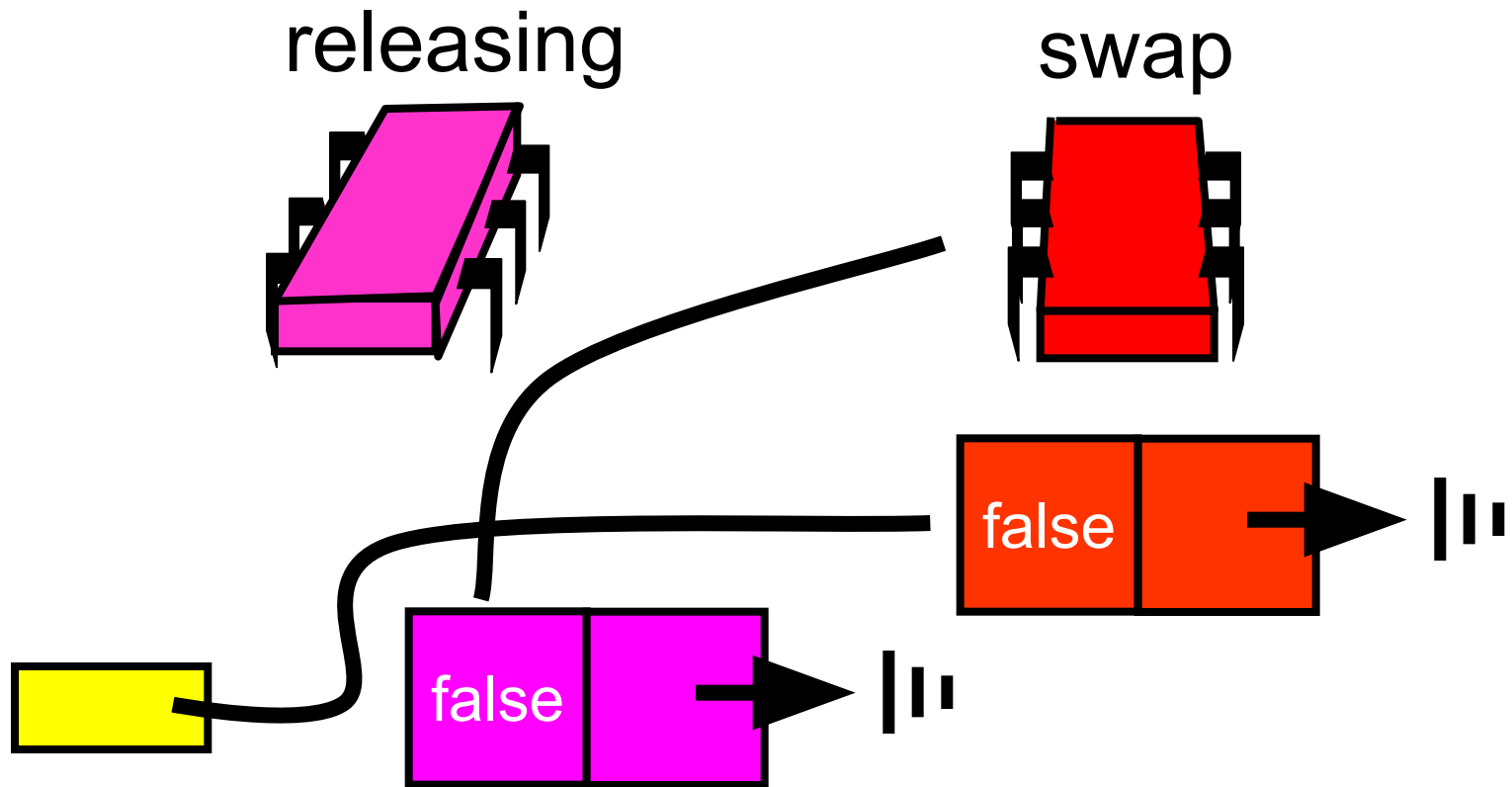
# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void lock() {
  QNode qnode = new QNode();
  QNode pred = tail.getAndSet(qnode);
  if (pred != null) {
   qnode.locked = true;
   pred.next = qnode;
   while (qnode.locked) {}
  }}}
```
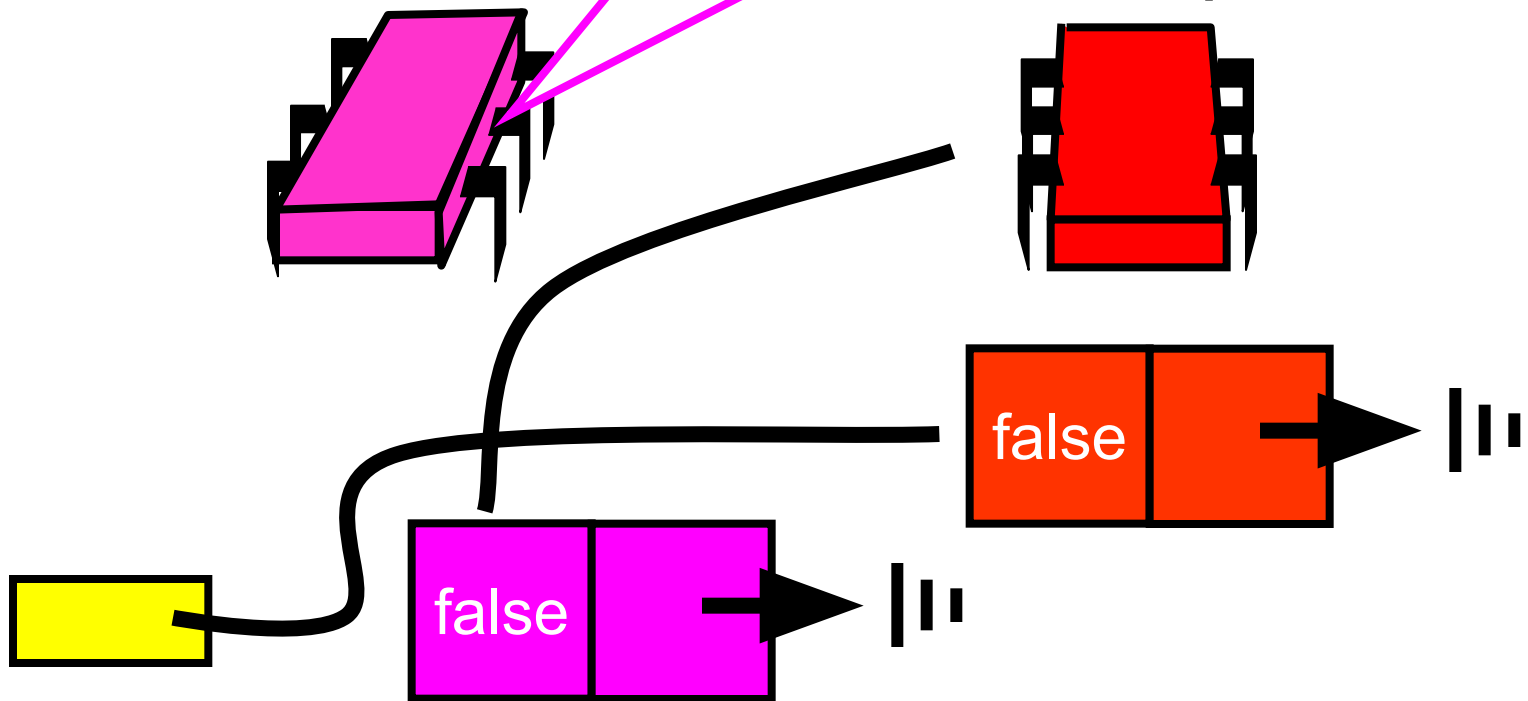
# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void lock() {
  QNode qnode = new QNode();
  QNode pred = tail.getAndSet(qnode);
  if (pred != null) {
   qnode.locked = true;
   pred.next = qnode;
   while (qnode.locked) {}
  }}}
```

**Make a QNode**

# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void lock() {
  QNode qnode = new QNode();
  QNode pred = tail.getAndSet(qnode);
  if (pred != null) {
   qnode.locked = true;
   pred.next = qnode;
   while (qnode.locked) {}
}}}
```

**add my Node to the tail of queue**

# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void lock() {
  QNode qnode = new QNode();
  QNode pred = tail.getAndSet(qnode);
  if (pred != null) {
   qnode.locked = true;
   pred.next = qnode;
   while (qnode.locked) {}
}}}
```

**Fix if queue was non-empty**

# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void lock() {
  QNode qnode = new QNode();
  QNode pred = tail.getAndSet(qnode);
  if (pred != null) {
   qnode.locked = true;
   pred.next = qnode;
   while (qnode.locked) {}
}}}
```

**Wait until unlocked**

# Purple Release

releasing        swap

false

false

# Purple Release

releasing          prepare to spin



true

false

# Purple Release

releasing         spinning

true

false

# Purple Release

releasing
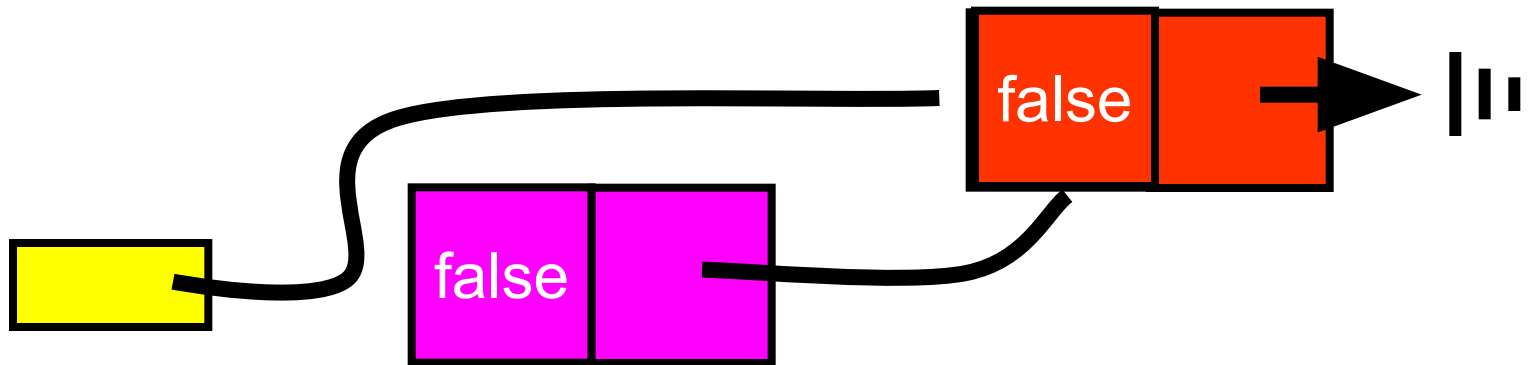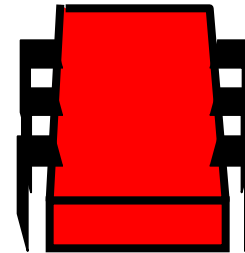
spinning

false

false

# Purple Release

releasing      Acquired lock

false

false

# MCS Queue Unlock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void unlock() {
   if (qnode.next == null) {
    if (tail.CAS(qnode, null)
     return;
    while (qnode.next == null) {}
   }
 qnode.next.locked = false;
}}
```
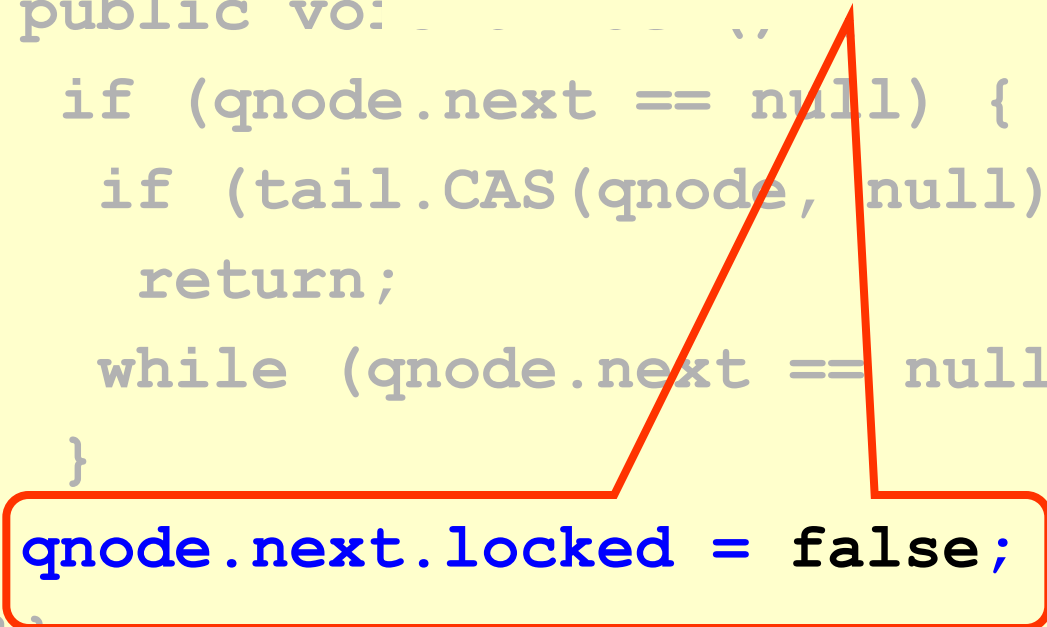
# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicReference tail;
 public void unlock() {
   if (qnode.next == null) {
     if (tail.CAS(qnode, null)
       return;
     while (qnode.next == null) {}
   }
 qnode.next.locked = false;
}}
```

**Missing successor?**

# MCS Queue Lock

**If really no successor, return**

```
                                    :k {

  if (qnode.next == null) {

    if (tail.CAS(qnode, null)

    return;

   while (qnode.next == null) {}
  }
 qnode.next.locked = false;
}}
```

# MCS Queue Lock

**Otherwise wait for successor to catch up**

```
                                            k {
  public void unlock() {
    if (qnode.next == null) {
      if (tail.CAS(qnode, null)
        return;
      while (qnode.next == null) {}
    }
    qnode.next.locked = false;
}}
```

# MCS Queue Lock

```
class MCSLock implements Lock {
 AtomicRef
 public vo
   if (qnode.next == null) {
     if (tail.CAS(qnode, null)
       return;
     while (qnode.next == null) {}
   }
   qnode.next.locked = false;
}}
```

**Pass lock to successor**

# Abortable Locks

- What if you want to give up waiting for a lock?

- For example
  - Timeout
  - Database transaction aborted by user

# Back-off Lock

- Aborting is trivial
  - Just return from lock() call
- Extra benefit:
  - No cleaning up
  - Wait-free
  - Immediate return

# Queue Locks

- Can't just quit
  - Thread in line behind will starve
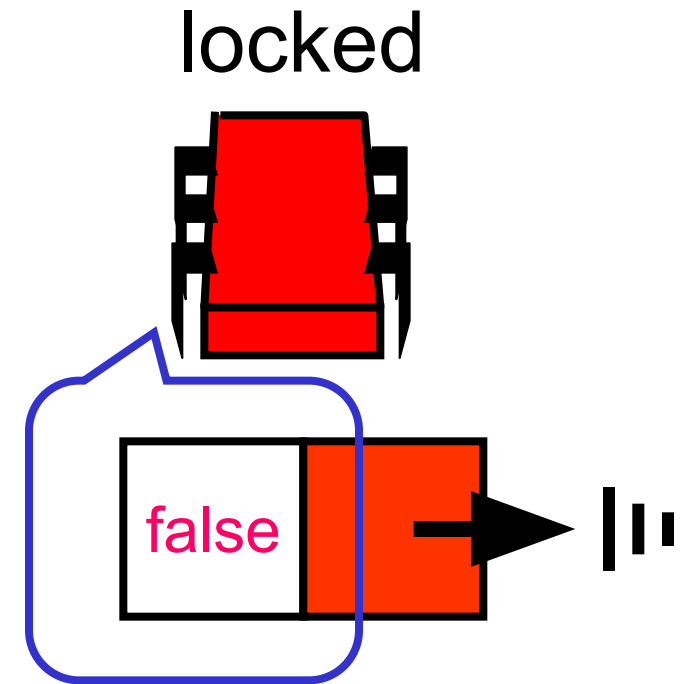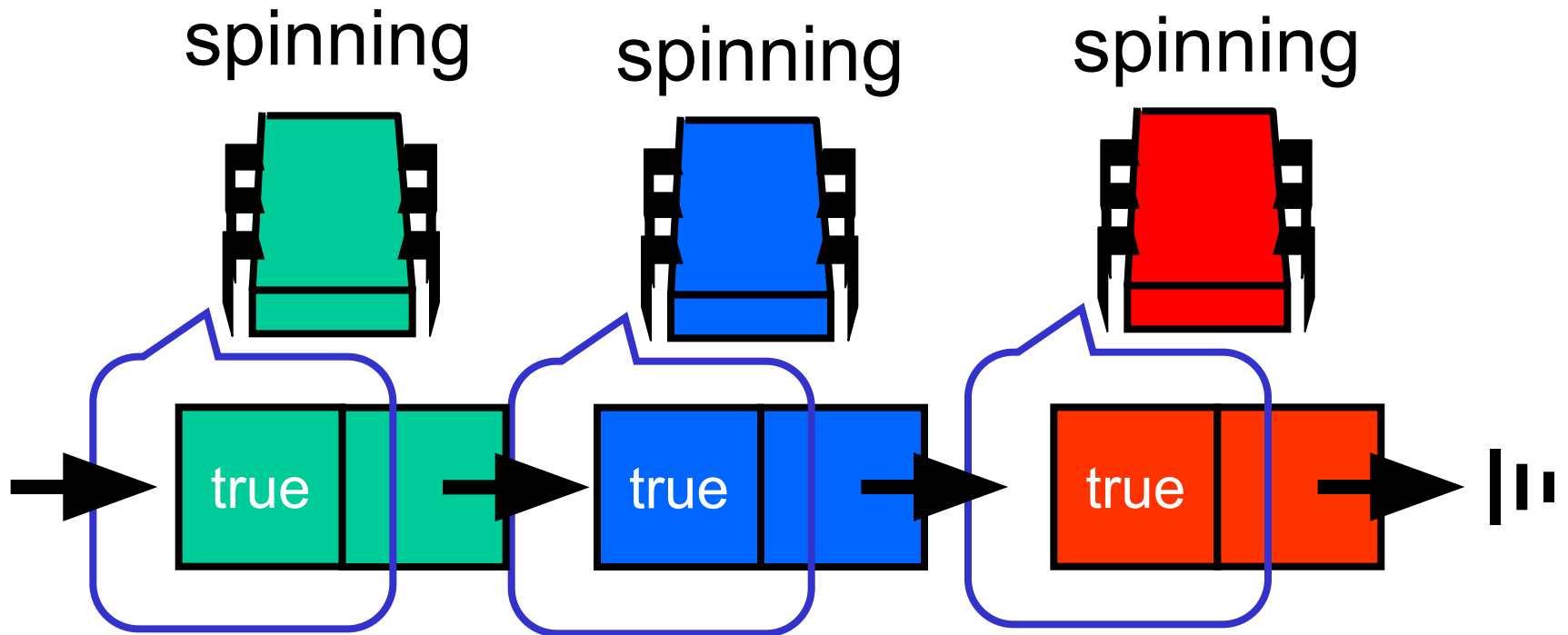- Need a graceful way out

# Queue Locks

# Queue Locks



locked     spinning     spinning

false     true     true

# Queue Locks

locked

spinning

false

true

# Queue Locks

locked

false

# Queue Locks

spinning        spinning        spinning

# Queue Locks

spinning

spinning

true    true    true

# Queue Locks

locked                    spinning



false → true → true → ‖▪

# Queue Locks

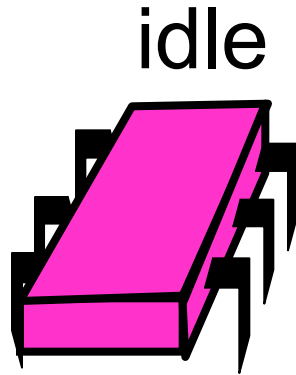spinning

false | true

# Queue Locks

pwned

false    true

# Abortable CLH Lock

- When a thread gives up
  - Removing node in a wait-free way is hard
- Idea:
  - let successor deal with it.

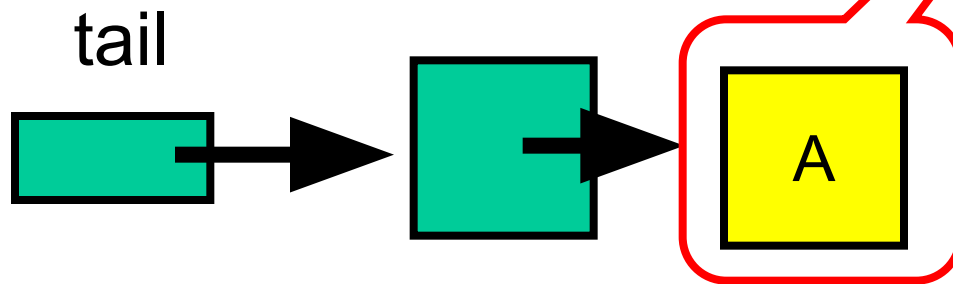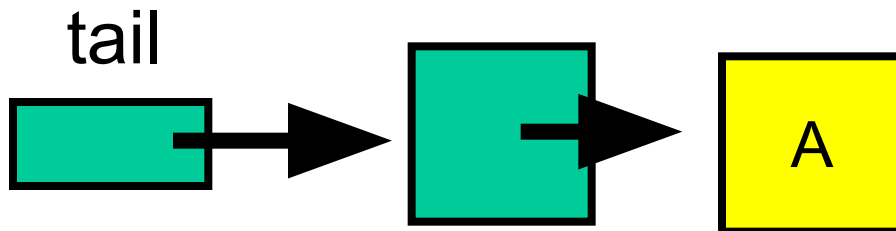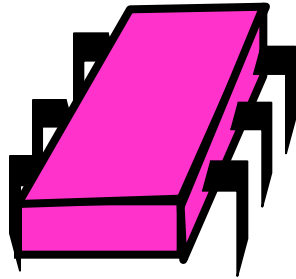# Initially

idle



Pointer to predecessor (or null)

tail

A

# Initially

idle



Distinguished available node means lock is free

tail



A

# Acquiring

acquiring
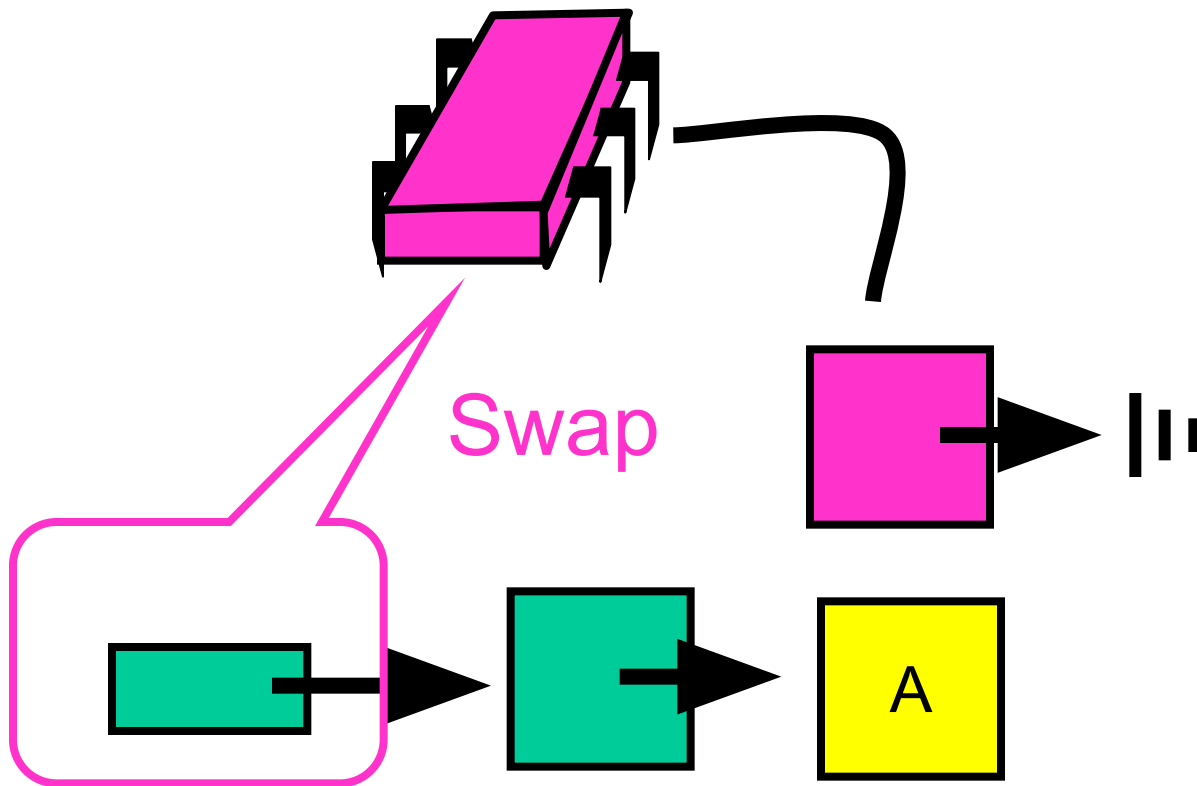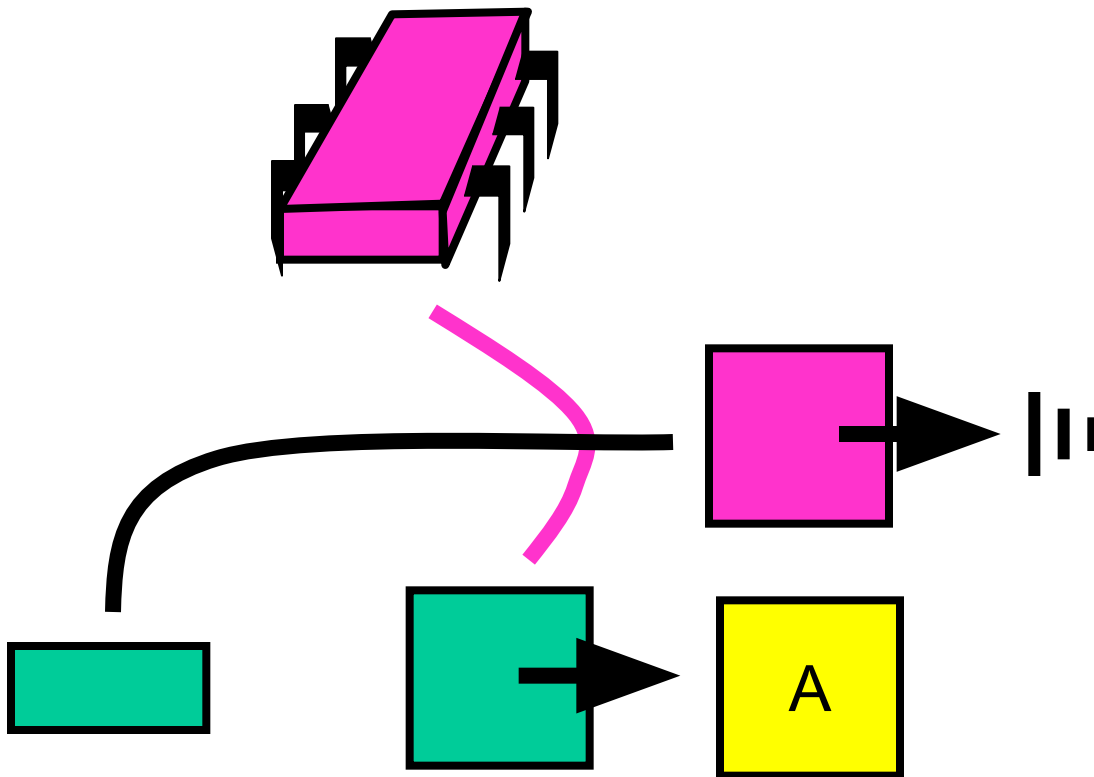


tail



A

# Acquiring

# Acquiring

acquiring

Swap

A

# Acquiring

acquiring

# Acquired

locked

Reference to **AVAILABLE** means **lock is free**.

A

# Normal Case



locked      spinning      spinning

**Null** means lock is
not free & request
not aborted

# One Thread Aborts

locked          Timed out          spinning

# Successor Notices

locked    Timed out    spinning

**Non-Null means predecessor aborted**

# Recycle Predecessor's Node

locked

spinning

# Spin on Earlier Node

locked

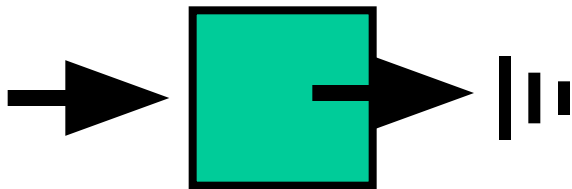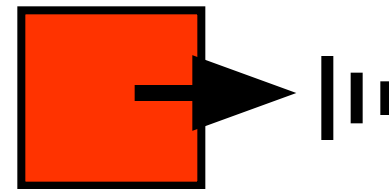spinning

# Spin on Earlier Node

released

spinning

A

*The lock is now mine*

# Time-out Lock

```
public class TOLock implements Lock {
   static QNode AVAILABLE
      = new QNode();
   AtomicReference<QNode> tail;
   ThreadLocal<QNode> myNode;
```

# Time-out Lock

```
public class TOLock implements Lock {
  static QNode AVAILABLE
    = new QNode();
  AtomicReference<QNode> tail;
  ThreadLocal<QNode> myNode;
```

**AVAILABLE node
signifies free lock**

# Time-out Lock

```
public class TOLock implements Lock {
  static QNode AVAILABLE
    = new QNode();
  AtomicReference<QNode> tail;
  ThreadLocal<QNode> myNode;
```

**Tail of the queue**

# Time-out Lock

```
public class TOLock implements Lock {
  static QNode AVAILABLE
    = new QNode();
  AtomicReference<QNode> tail;
  ThreadLocal<QNode> myNode;
```

**Remember my node …**

# Time-out Lock

```
public boolean lock(long timeout) {
    QNode qnode = new QNode();
    myNode.set(qnode);
    qnode.prev = null;
    QNode myPred = tail.getAndSet(qnode);
    if (myPred== null
          || myPred.prev == AVAILABLE) {
        return true;
    }
…
```

# Time-out Lock

```
public boolean lock(long timeout) {
  QNode qnode = new QNode();
  myNode.set(qnode);
  qnode.prev = null;
  QNode myPred = tail.getAndSet(qnode);
  if (myPred == null
       || myPred.prev == AVAILABLE) {
      return true;
    }
```

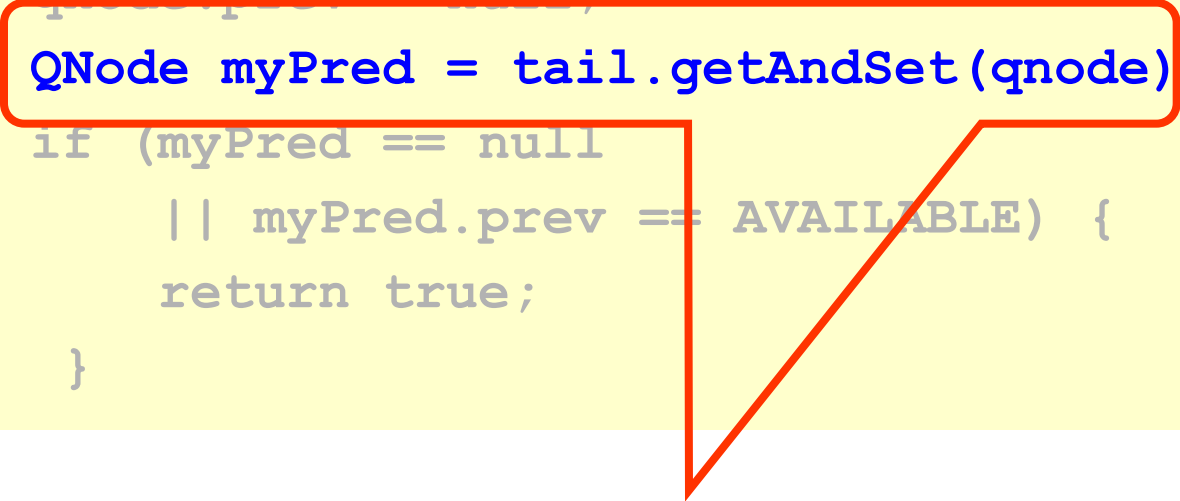**Create & initialize node**

# Time-out Lock

```
public boolean lock(long timeout) {
  QNode qnode = new QNode();
  myNode.set(qnode);
  qnode.prev = null;
  QNode myPred = tail.getAndSet(qnode);
  if (myPred == null
       || myPred.prev == AVAILABLE) {
      return true;
  }
```
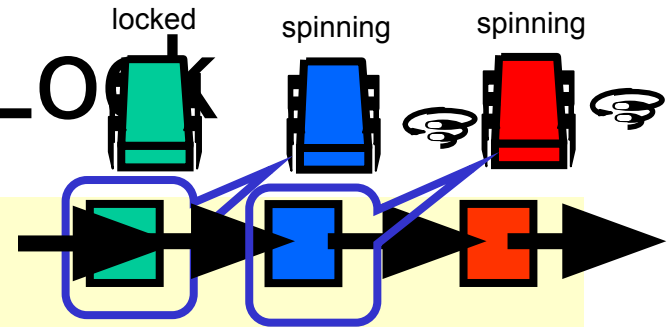
**Swap with tail**

# Time-out Lock

```
public boolean lock(long timeout) {
    QNode qnode = new QNode();
    myNode.set(qnode);
    qnode.prev = null;
    QNode myPred = tail.getAndSet(qnode);
    if (myPred == null
          || myPred.prev == AVAILABLE) {
        return true;
    }
    ...
```

**If predecessor absent or released, we are done**

# Time-out Lock

```
…
    long start = now();
    while (now()- start < timeout) {
        QNode predPred = myPred.prev;
        if (predPred == AVAILABLE) {
            return true;
        } else if (predPred != null) {
            myPred = predPred;
        }
    }
    …
```
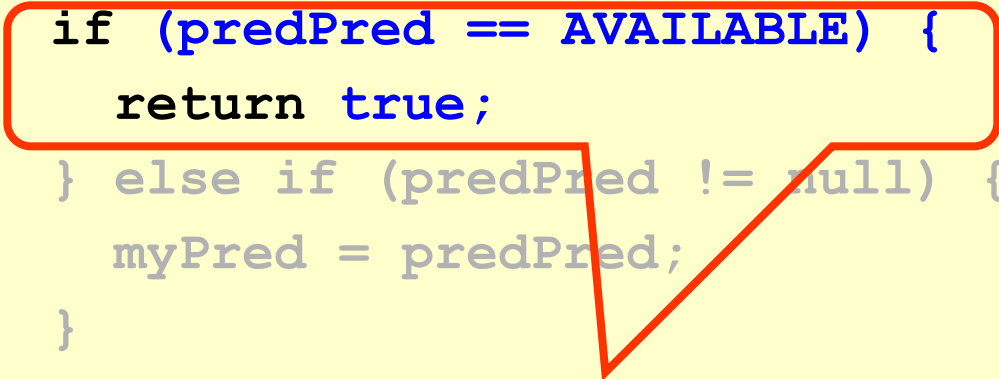
# Time-out Lock

```
…
long start = now();
while (now()- start < timeout) {
  QNode predPred = myPred.prev;
  if (predPred == AVAILABLE) {
    return true;
  } else if (predPred != null) {
    myPred = predPred;
  }
}
…
```

**Keep trying for a while**

**…**

# Time-out Lock

```
…
   long start = now();
   while (now()- start < timeout) {
      QNode predPred = myPred.prev;
      if (predPred == AVAILABLE) {
         return true;
      } else if (predPred != null) {
        myPred = predPred;
      }
   }
   …
```

**Spin on predecessor's prev field**

# Time-out Lock

```
…
  long start = now();
  while (now()- start < timeout) {
    QNode predPred = myPred.prev;
    if (predPred == AVAILABLE) {
      return true;
    } else if (predPred != null) {
      myPred = predPred;
    }
  }
…
```

**Predecessor released lock**

# Time-out Lock

```
…
  long start = now();
  while (now()- start < timeout) {
    QNode predPred = myPred.prev;
    if (predPred == AVAILABLE) {
      return true;
    } else if (predPred != null) {
      myPred = predPred;
    }
  }
  …
```

**Predecessor aborted, advance one**

# Time-out Lock

```
…
if (!tail.compareAndSet(qnode, myPred))
    qnode.prev = myPred;
    return false;
  }
}
```
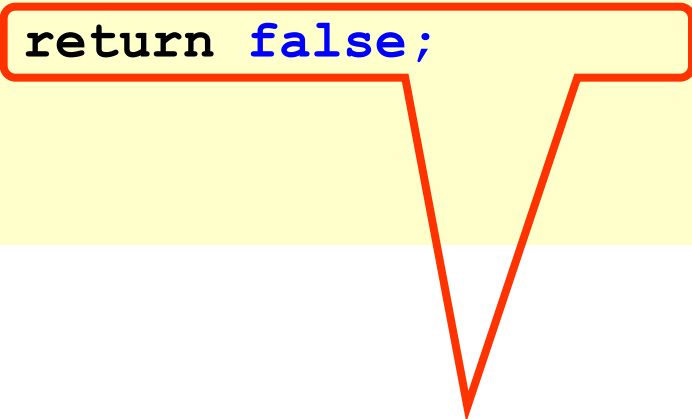
## What do I do when I time out?

# Time-out Lock

```
…
if (!tail.compareAndSet(qnode, myPred))
    qnode.prev = myPred;
    return false;
  }
}
```

**Do I have a successor?**
**If CAS fails, I do.**
**Tell it about myPred**

# Time-out Lock

```
…
if (!tail.compareAndSet(qnode, myPred))
    qnode.prev = myPred;
    return false;
    }
}
```

**If CAS succeeds: no successor, simply return false**

# Time-Out Unlock

```
public void unlock() {
  QNode qnode = myNode.get();
  if (!tail.compareAndSet(qnode, null))
    qnode.prev = AVAILABLE;
}
```
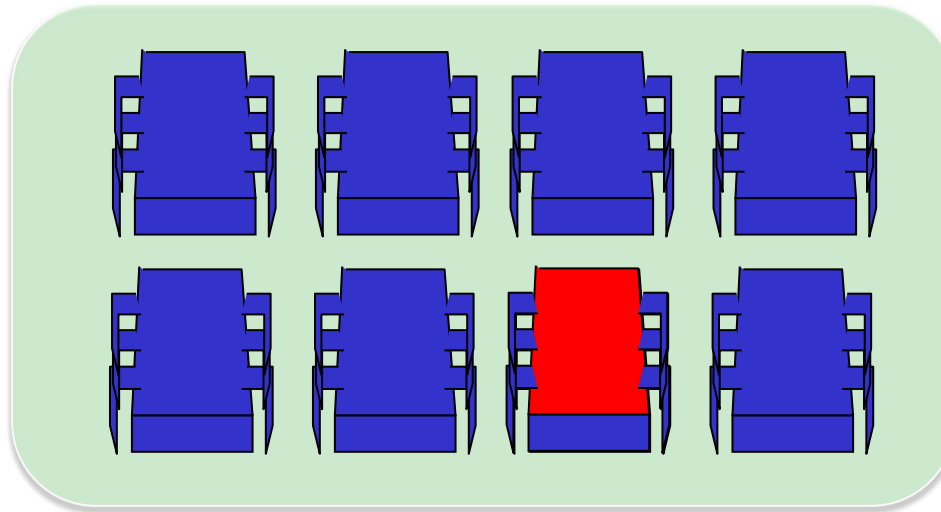
# Time-out Unlock

```
public void unlock() {
  QNode qnode = myNode.get();
  if (!tail.compareAndSet(qnode, null))
    qnode.prev = AVAILABLE;
}
```

**If CAS failed: successor exists, notify it can enter**

# Timing-out Lock

```
public void unlock() {
   QNode qnode = myNode.get();
   if (!tail.compareAndSet(qnode, null))
      qnode.prev = AVAILABLE;
}
```

**CAS successful: set tail to null, no clean up since no successor waiting**

# Fairness and NUMA Locks

- MCS lock mechanics are aware of NUMA

- Lock Fairness is FCFS

- Is this a good fit with NUMA and Cache-Coherent NUMA machines?

# Lock Data Access in NUMA Machine

# "Who's the Unfairest of Them All?"



- locality crucial to NUMA performance
- Big gains if threads from same node/cluster obtain lock consecutively
- Unfairness pays

# Hierarchical Backoff Lock (HBO)



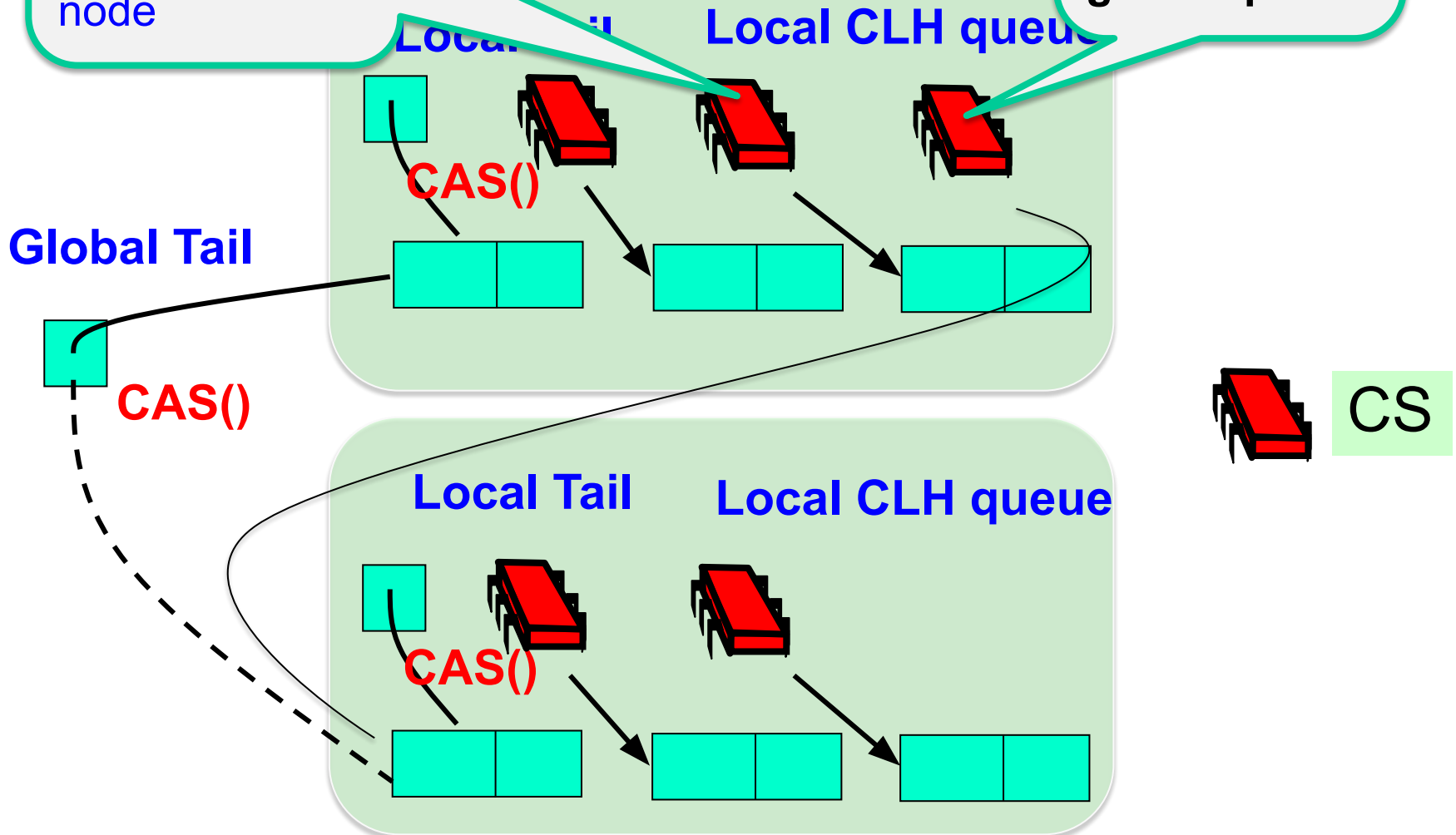Back off less for thread from same node

Unfairness is key to performance

time
4d    2d    d

time
4d    2d    d

CS

Global T&T&S lock

# Hierarchical Backoff Lock (HBO)

- Advantages:
  - Simple, improves locality
- Disadvantages:
  - Requires platform specific tuning
  - Unstable
  - Unfair
  - Continuous invalidations on shared global lock word

# chical CLH Lock (

Each thread spins on cached copy of predecessor's node

Thread at local head splices **local queue into global queue**

**Local Tail**

**Local CLH queue**

CAS()

**Global Tail**

CAS()

CS

**Local Tail**

**Local CLH queue**

CAS()

# Hierarchical CLH Lock (HCLH)



**Threads access 4 cache lines in CS**

# Hierarchical CLH Lock (HCLH)

- Advantages:
  - Improved locality
  - Local spinning
  - Fair
- Disadvantages:
  - Complex code implies long common path
  - Splicing into both local and global requires CAS
  - Hard to get long local sequences

"Nothing yet. ... How about you, Newton?"

# Lock Cohorting

- General technique for converting almost any lock into a NUMA lock

- Allows combining different lock types

- But need these locks to have certain properties (will discuss shortly)

# Lock Cohorting

# Thread Obliviousness

- A lock is ***thread-oblivious*** if
  - After being acquired by one thread,
  - Can be released by another

# Cohort Detection

- A lock x provides **_cohort detection_** if
  - It can tell whether any thread is trying to acquire it

# Lock Cohorting

- Two levels of locking
- **Global lock**: thread oblivious
  - Thread acquiring the lock can be different than one releasing it
- **Local lock**: cohort detection
  - Thread releasing can detect if some thread is waiting to acquire it

Two new states: *acquire local* and *acquire global*. ?Do we own global lock

In MCS Lock, *cohort detection* by checking successor pointer

**Local MCS lock tail**

**CAS()**

False

False ← True

Bound number of consecutive acquires to control unfairness

**Loca**

**CAS()**

False ← False ← True

Global backoff lock

CS

time **4d** **2d d**

BO Lock is *thread oblivious* by definition

# Lock Co... BO Lock



How to add *cohort detection property?* to BO lock

4d    2d d

4d    2d d

Global backoff lock

CS

time **4d**    **2d**

As noted BO Lock is *thread oblivious*

# Lock Co... -BO Lock

Add successorExists field before attempting to acquire local lock.

successorExists reset on lock release.

Release might overwrite another successor's write … but we don't care…why?



4d   2d d

4d   2d d

time  4d   2d d

CS

# C-B...e-Out



Aborting thread resets successorExists field before leaving local lock. Spinning threads set it to true.
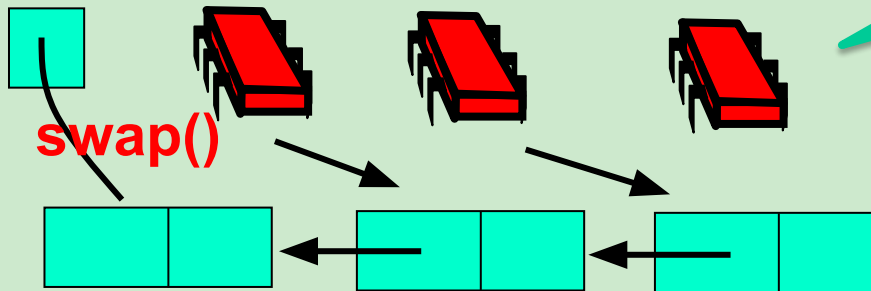
BO locks trivially abortable

Global backoff lock

CS

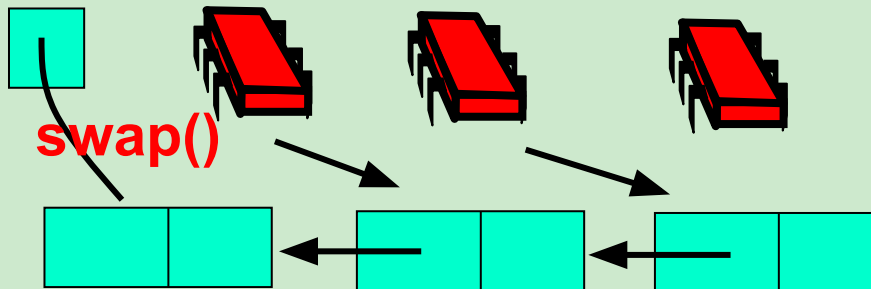If releasing thread finds successorExists false, it releases global lock

4d    2d d

4d    2d d

2d d

# Time-Out NUMA Lock



Local time-out queue

swap()

Local time-out queue

swap()

Local Time-Out locks have cohort detection ?property …why …Not enough

CS

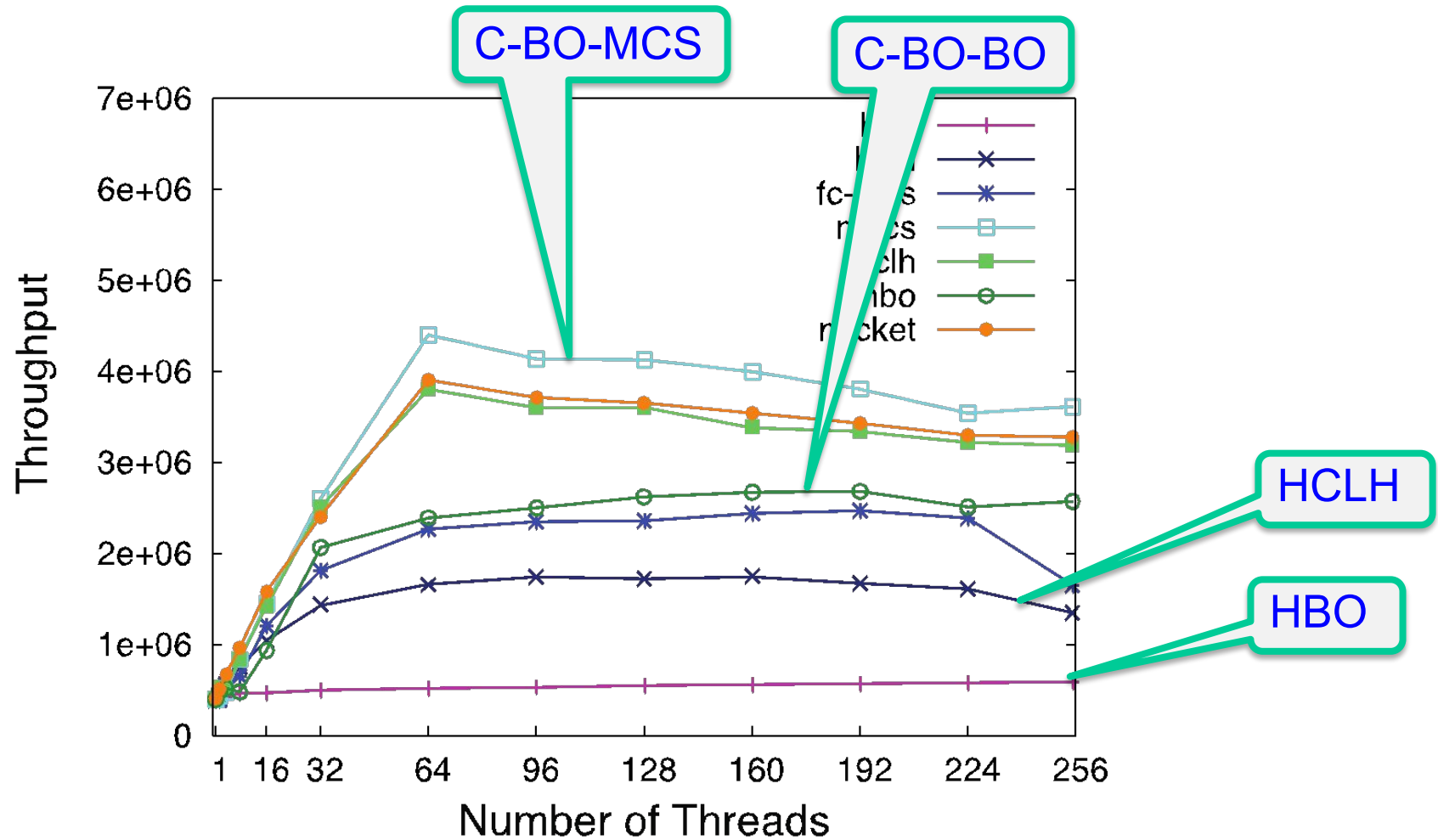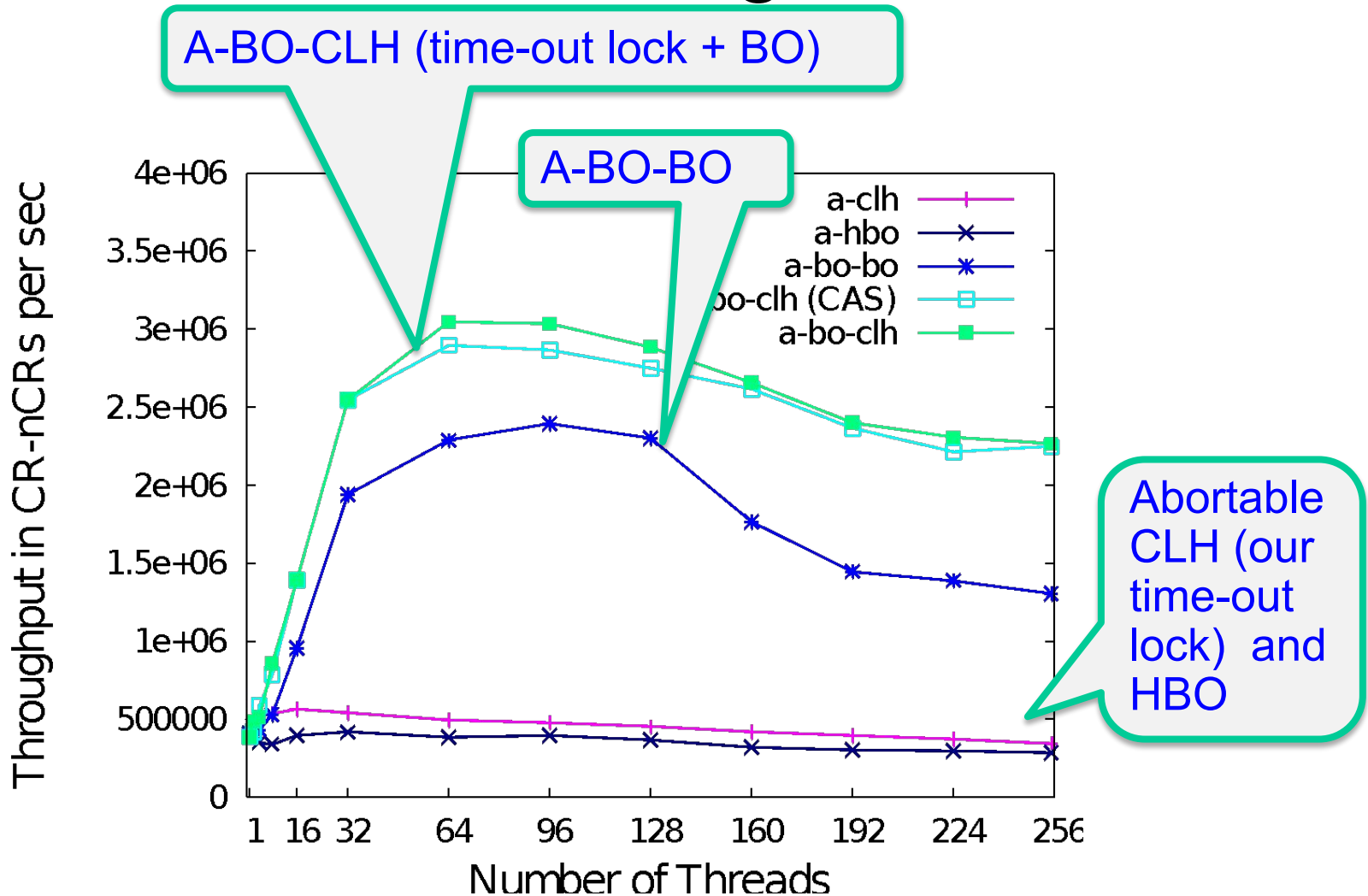time **4d** **2d d**

# Lock Cohorting

- Advantages:
  - Great locality
  - Low contention on shared lock
  - Practically no tuning
  - Has whatever properties you want:
    - Can be more or less fair, abortable…
    just choose the appropriate type of locks…
- Disadvantages:
  - Must tune fairness parameters

# Lock Cohorting

# Time-Out (Abortable) Lock Cohorting

# One Lock To Rule Them All?

- TTAS+Backoff, CLH, MCS, ToLock…
- Each better than others in some way
- There is no one solution
- Lock we pick really depends on:
  - the application
  - the hardware
  - which properties are important