

Design of Parallel and High-Performance Computing

Fall 2014

Lecture: Locks and Lock-Free continued

Instructor: Torsten Hoefler & Markus Püschel

TA: Timo Schneider, Arnamoy Bhattacharyya

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Administrivia

- **Intermediate (very short) presentation: Thursday 11/27 during recitation**
 - Should have first results and a real plan!
 - Time to get very brief feedback
 - Focus on:
 - What tools/programming language/parallelization scheme do you use?*
 - Which architecture? (we only offer access to Xeon Phi, you may use different)*
 - How to verify correctness of the parallelization?*
 - How to argue about performance (bounds, what to compare to?)*
 - (Somewhat) realistic use-cases and input sets?*
 - What are the key concepts employed?*
 - What are the main obstacles?*
- **Final project presentation: Monday 12/15 during last lecture**
 - Report will be due in January!
 - *Still, starting to write early is very helpful -- write – rewrite – rewrite ...*

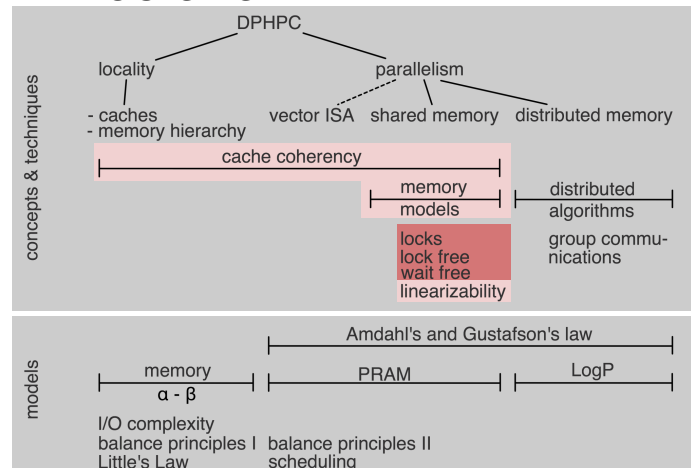
2

Review of last lecture

- **Hardware operations for concurrency control**
- **More on locks (using advanced operations)**
 - Spin locks
 - Various optimized locks
- **Stopped right before queue locks (the best ones 😊)**

3

DPHPC Overview



4

Goals of this lecture

- **Even more on locks (issues and extended concepts)**
 - Queue locks
 - Deadlocks, priority inversion, competitive spinning, semaphores
- **Case studies**
 - Barrier
- **Locks in practice: a set structure**

5

Improvements?

- **Are TAS locks perfect?**
 - What are the two biggest issues?
 - Cache coherency traffic (contending on same location with expensive atomics)
 - or --
 - Critical section underutilization (waiting for backoff times will delay entry to CR)
- **What would be a fix for that?**
 - How is this solved at airports and shops (often at least)?
- **Queue locks -- Threads enqueue**
 - Learn from predecessor if it's their turn
 - Each threads spins at a different location
 - FIFO fairness

6

Array Queue Lock

■ Array to implement queue

- Tail-pointer shows next free queue position
- Each thread spins on own location
CL padding!
- `index[]` array can be put in TLS

■ So are we done now?

- What's wrong?
- Synchronizing M objects requires $\Theta(NM)$ storage
- What do we do now?

```
volatile int array[n] = {1,0,...,0};
volatile int index[n] = {0,0,...,0};
volatile int tail = 0;

void lock() {
    index[tid] = GetAndInc(tail) % n;
    while (!array[index[tid]]); // wait to receive lock
}

void unlock() {
    array[index[tid]] = 0; // I release my lock
    array[(index[tid] + 1) % n] = 1; // next one
}
```

7

CLH Lock (1993)

■ List-based (same queue principle)

■ Discovered twice by Craig, Landin, Hagersten 1993/94

■ 2N+3M words

- N threads, M locks

■ Requires thread-local qnode pointer

- Can be hidden!

```
typedef struct qnode {
    struct qnode *prev;
    int succ_blocked;
} qnode;

qnode *lck = new qnode; // node owned by lock

void lock(qnode *lck, qnode *qn) {
    qn->succ_blocked = 1;
    qn->prev = FetchAndSet(lck, qn);
    while (qn->prev->succ_blocked);
}

void unlock(qnode **qn) {
    qnode *pred = (*qn)->prev;
    (*qn)->succ_blocked = 0;
    *qn = pred;
}
```

8

CLH Lock (1993)

■ Qnode objects represent thread state!

- `succ_blocked == 1` if waiting or acquired lock
- `succ_blocked == 0` if released lock

■ List is implicit!

- One node per thread
- Spin location changes
NUMA issues (cacheless)

■ Can we do better?

```
typedef struct qnode {
    struct qnode *prev;
    int succ_blocked;
} qnode;

qnode *lck = new qnode; // node owned by lock

void lock(qnode *lck, qnode *qn) {
    qn->succ_blocked = 1;
    qn->prev = FetchAndSet(lck, qn);
    while (qn->prev->succ_blocked);
}

void unlock(qnode **qn) {
    qnode *pred = (*qn)->prev;
    (*qn)->succ_blocked = 0;
    *qn = pred;
}
```

9

MCS Lock (1991)

■ Make queue explicit

- Acquire lock by appending to queue
- Spin on own node until locked is reset

■ Similar advantages as CLH but

- Only $2N + M$ words
- Spinning position is fixed!
Benefits cache-less NUMA

■ What are the issues?

- Releasing lock spins
- More atomics!

```
typedef struct qnode {
    struct qnode *next;
    int succ_blocked;
} qnode;

qnode *lck = NULL;

void lock(qnode *lck, qnode *qn) {
    qn->next = NULL;
    qnode *pred = FetchAndSet(lck, qn);
    if(pred != NULL) {
        qn->locked = 1;
        pred->next = qn;
        while(qn->locked);
    }
}

void unlock(qnode *lck, qnode *qn) {
    if(qn->next == NULL) { // if we're the last waiter
        if(CAS(lck, qn, NULL)) return;
        while(qn->next == NULL); // wait for pred arrival
    }
    qn->next->locked = 0; // free next waiter
    qn->next = NULL;
}
```

Lessons Learned!

■ Key Lesson:

- Reducing memory (coherency) traffic is most important!
- Not always straight-forward (need to reason about CL states)

■ MCS: 2006 Dijkstra Prize in distributed computing

- *"an outstanding paper on the principles of distributed computing, whose significance and impact on the theory and/or practice of distributed computing has been evident for at least a decade"*
- *"probably the most influential practical mutual exclusion algorithm ever"*
- *"vastly superior to all previous mutual exclusion algorithms"*
- fast, fair, scalable → widely used, always compared against!

11

Time to Declare Victory?

■ Down to memory complexity of $2N+M$

- Probably close to optimal (try to prove a lower bound!)

■ Only local spinning

- Several variants with low expected contention

■ But: we assumed sequential consistency ☹

- Reality causes trouble sometimes
- Sprinkling memory fences may harm performance
- Open research on minimally-synching algorithms!
Come and talk to me if you're interested

12

More Practical Optimizations

- Let's step back to "data race"
 - (recap) two operations A and B on the same memory cause a data race if one of them is a write ("conflicting access") and neither $A \rightarrow B$ nor $B \rightarrow A$
 - So we put conflicting accesses into a CR and lock it!
This also guarantees memory consistency in C++/Java!
- Let's say you implement a web-based encyclopedia
 - Consider the "average two accesses" – do they conflict?

Reader-Writer Locks

- Allows multiple concurrent reads
 - Multiple reader locks concurrently in CR
 - Guarantees mutual exclusion between writer and writer locks and reader and writer locks
- Syntax:
 - `read_(un)lock()`
 - `write_(un)lock()`

A Simple RW Lock

- Seems efficient!?
 - Is it? What's wrong?
 - Polling CAS!
- Is it fair?
 - Readers are preferred!
 - Can always delay writers (again and again and again)

```
const W = 1;
const R = 2;
volatile int lock=0; // LSB is writer flag!

void read_lock(lock_t lock) {
  AtomicAdd(lock, R);
  while(lock & W);
}

void write_lock(lock_t lock) {
  while(!CAS(lock, 0, W));
}

void read_unlock(lock_t lock) {
  AtomicAdd(lock, -R);
}

void write_unlock(lock_t lock) {
  AtomicAdd(lock, -W);
}
```

Fixing those Issues?

- Polling issue:
 - Combine with MCS lock idea of queue polling
- Fairness:
 - Count readers and writers

(1991) Scalable Reader-Writer Synchronization for Shared-Memory Multiprocessors
John M. Mellor-Crummey¹, Michael L. Scott²
¹IBM Research Division, Almaden Research Center, San Jose, CA 95120
²Computer Science Department, University of Rochester, Rochester, NY 14627

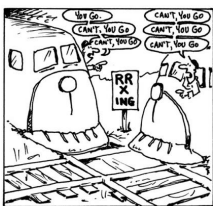
Abstract
Reader-writer synchronization relaxes the constraints of mutual exclusion to permit more than one process to inspect a shared object concurrently, using no more than a single compare-and-swap (CAS) instruction. In this paper, we describe a new reader-writer lock that is designed to be scalable to high processor counts. However, we describe multiprocessors as other advantages to have process keep wait. Unfortunately, implementation of these two locks on shared-memory multiprocessors typically cause memory and lock contention that degrades performance. Several researchers have shown how to improve scalable mutual exclusion locks that exploit locality in the memory hierarchy of shared-memory multiprocessors to eliminate contention for memory and for the processor-memory interface. In this paper, we present reader-writer locks that naturally exploit locality to achieve scalability, with variable reader performance, writer performance, and reader-writer mutual exclusion. Performance results on a 32-processor multiprocessor demonstrate that our algorithms provide low latency and excellent scalability.

1. Introduction
Reader-writer synchronization relaxes the constraints of mutual exclusion to permit more than one process to inspect a shared object concurrently, using no more than a single compare-and-swap (CAS) instruction. In this paper, we describe a new reader-writer lock that is designed to be scalable to high processor counts. However, we describe multiprocessors as other advantages to have process keep wait. Unfortunately, implementation of these two locks on shared-memory multiprocessors typically cause memory and lock contention that degrades performance. Several researchers have shown how to improve scalable mutual exclusion locks that exploit locality in the memory hierarchy of shared-memory multiprocessors to eliminate contention for memory and for the processor-memory interface. In this paper, we present reader-writer locks that naturally exploit locality to achieve scalability, with variable reader performance, writer performance, and reader-writer mutual exclusion. Performance results on a 32-processor multiprocessor demonstrate that our algorithms provide low latency and excellent scalability.

The final algorithm (Alg. 4) has a flaw that was corrected in 2003!

Deadlocks

- Kansas state legislature: "When two trains approach each other at a crossing, both shall come to a full stop and neither shall start up again until the other has gone."
[according to Botkin, Harlow "A Treasury of Railroad Folklore" (pp. 381)]



What are necessary conditions for deadlock?

Deadlocks

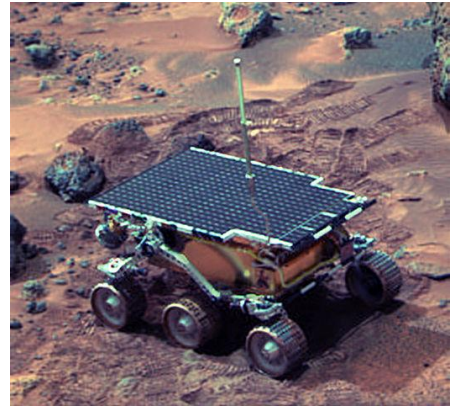
- Necessary conditions:
 - Mutual Exclusion
 - Hold one resource, request another
 - No preemption
 - Circular wait in dependency graph
- One condition missing will prevent deadlocks!
 - Different avoidance strategies (which?)

Issues with Spinlocks

- **Spin-locking is very wasteful**
 - The spinning thread occupies resources
 - Potentially the PE where the waiting thread wants to run → requires context switch!
- **Context switches due to**
 - Expiration of time-slices (forced)
 - Yielding the CPU

19

What is this?



20

Why is the 1997 Mars Rover in our lecture?

- **It landed, received program, and worked ... until it spuriously rebooted!**
 - → watchdog
- **Scenario (vxWorks RT OS):**
 - Single CPU
 - Two threads A,B sharing common bus, using locks
 - (independent) thread C wrote data to flash
 - Priority: A→C→B (A highest, B lowest)
 - Thread C would run into a livelock (infinite loop)
 - Thread B was preempted by C while holding lock
 - Thread A got stuck at lock ☹

[http://research.microsoft.com/en-us/um/people/mbj/Mars_Pathfinder/Authoritative_Account.html]

21

Priority Inversion

- **If busy-waiting thread has higher priority than thread holding lock ⇒ no progress!**
- **Can be fixed with the help of the OS**
 - E.g., mutex priority inheritance (temporarily boost priority of task in CR to highest priority among waiting tasks)

22

Condition Variables

- **Allow threads to yield CPU and leave the OS run queue**
 - Other threads can get them back on the queue!
- **cond_wait(cond, lock) – yield and go to sleep**
- **cond_signal(cond) – wake up sleeping threads**
- **Wait and signal are OS calls**
 - Often expensive, which one is more expensive?
Wait, because it has to perform a full context switch

23

Condition Variable Semantics

- **Hoare-style:**
 - Signaler passes lock to waiter, signaler suspended
 - Waiter runs immediately
 - Waiter passes lock back to signaler if it leaves critical section or if it waits again
- **Mesa-style (most used):**
 - Signaler keeps lock
 - Waiter simply put on run queue
 - Needs to acquire lock, may wait again

24

When to Spin and When to Block?

- **Spinning consumes CPU cycles but is cheap**
 - “Steals” CPU from other threads
- **Blocking has high one-time cost and is then free**
 - Often hundreds of cycles (trap, save TCB ...)
 - Wakeup is also expensive (latency)
Also cache-pollution
- **Strategy:**
 - Poll for a while and then block

25

When to Spin and When to Block?

- **What is a “while”?**
- **Optimal time depends on the future**
 - When will the active thread leave the CR?
 - Can compute optimal offline schedule
 - Actual problem is an online problem
- **Competitive algorithms**
 - An algorithm is c -competitive if for a sequence of actions x and a constant a holds:
$$C(x) \leq c * C_{opt}(x) + a$$
 - What would a good spinning algorithm look like and what is the competitiveness?

26

Competitive Spinning

- **If T is the overhead to process a wait, then a locking algorithm that spins for time T before it blocks is 2-competitive!**
 - Karlin, Manasse, McGeoch, Owicki: “Competitive Randomized Algorithms for Non-Uniform Problems”, SODA 1989
- **If randomized algorithms are used, then $e/(e-1)$ -competitiveness (~ 1.58) can be achieved**

27

Generalized Locks: Semaphores

- **Controlling access to more than one resource**
 - Described by Dijkstra 1965
- **Internal state is an atomic counter C**
- **Two operations:**
 - $P()$ – block until $C > 0$; decrement C (atomically)
 - $V()$ – signal and increment C
- **Binary or 0/1 semaphore equivalent to lock**
 - C is always 0 or 1, i.e., $V()$ will not increase it further
- **Trivia:**
 - If you’re lucky (aeheh, speak Dutch), mnemonics:
Verhogen (increment) and Prolaag (probeer te verlagen = try to reduce)

28

Semaphore Implementation

- **Can be implemented with mutual exclusion!**
 - And can be used to implement mutual exclusion ☺
- **... or with test and set and many others!**
- **Also has fairness concepts:**
 - Order of granting access to waiting (queued) threads
 - strictly fair (starvation impossible, e.g., FIFO)
 - weakly fair (starvation possible, e.g., random)

29

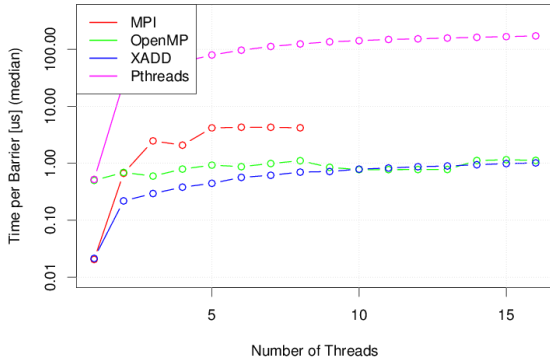
Case Study 1: Barrier

- **Barrier semantics:**
 - No process proceeds before all processes reached barrier
 - Similar to mutual exclusion but not exclusive, rather “synchronized”
- **Often needed in parallel high-performance programming**
 - Especially in SPMD programming style
- **Parallel programming “frameworks” offer barrier semantics (pthread, OpenMP, MPI)**
 - `MPI_Barrier()` (process-based)
 - `pthread_barrier`
 - `#pragma omp barrier`
 - `lock xadd + spin`
 - *Problem: when to re-use the counter?*
 - *Cannot just set it to 0 ☹*
 - *Trick: “lock xadd -1” when done ☺*

[cf. <http://www.spiral.net/software/barrier.html>]

30

Barrier Performance



31

Case Study 2: Reasoning about Semantics

Comments on a Problem in Concurrent Programming Control

Dear Editor:

I would like to comment on Mr. Dijkstra's solution [Solution of a problem in concurrent programming control. *Comm ACM* 8 (Sept. 1965), 569] to a messy problem that is hardly academic. We are using it now on a multiple computer complex.

When there are only two computers, the algorithm may be simplified to the following:

Boolean array $b(0; 1)$ **integer** $k, i, j,$

comment This is the program for computer i , which may be either 0 or 1, computer $j \neq i$ is the other one, 1 or 0;

$C0: b(i) := \text{false};$

$C1: \text{if } k \neq i \text{ then begin}$

$C2: \text{if not } b(j) \text{ then go to } C2;$

else $k := i;$ **go to** $C1$ **end;**

else critical section;

$b(i) := \text{true};$

remainder of program;

go to $C0;$

end

CACM
Volume 9 Issue 1, Jan. 1966

Mr. Dijkstra has come up with a clever solution to a really practical problem.

HARRIS HYMAN
Munttype
New York, New York

32

Case Study 2: Reasoning about Semantics

Is the proposed algorithm correct?

- We may prove it manually
 - Using tools from the last lecture
 - reason about the state space of H
- Or use automated proofs (model checking)
 - E.g., SPIN (Promela syntax)

```

bool want[2];
bool turn;
byte cnt;

proctype P(bool i)
{
    want[i] = 1;
    do
    :: (turn != i) ->
        (!want[1-i]);
        turn = i
    :: (turn == i) ->
        break
    od;
    skip; /* critical section */
    cnt = cnt+1;
    assert(cnt == 1);
    cnt = cnt-1;
    want[i] = 0
}

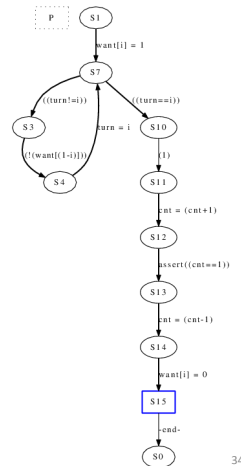
init { run P(0); run P(1) }
    
```

33

Case Study 2: Reasoning about Semantics

Spin tells us quickly that it found a problem

- A sequentially consistent order that violates mutual exclusion!
- It's not always that easy
 - This example comes from the SPIN tutorial
 - More than two threads make it much more demanding!
- More in the recitation!



34

Locks in Practice

Running example: List-based set of integers

- $S.insert(v)$ – return true if v was inserted
- $S.remove(v)$ – return true if v was removed
- $S.contains(v)$ – return true iff v in S

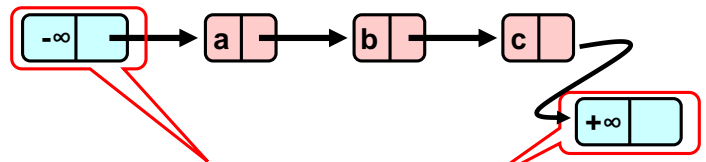
Simple ordered linked list

- Do not use this at home (poor performance)
- Good to demonstrate locking techniques
 - E.g., skip list would be faster but more complex

35

Set Structure in Memory

- This and many of the following illustrations are provided by Maurice Herlihy in conjunction with the book "The Art of Multiprocessor Programming"



Sorted with Sentinel nodes
(min & max possible keys)

36

Sequential Set

```
boolean add(S, x) {
    node *pred = S.head;
    node *curr = pred.next;
    while(curr.key < x) {
        pred = curr;
        curr = pred.next;
    }
    if(curr.key == x)
        return false;
    else {
        node n = new node();
        n.key = x;
        n.next = curr;
        pred.next = n;
    }
    return true;
}
```

```
boolean remove(S, x) {
    node *pred = S.head;
    node *curr = pred.next;
    while(curr.key < x) {
        pred = curr;
        curr = pred.next;
    }
    if(curr.key == x) {
        pred.next = curr.next;
        free(curr);
        return true;
    }
    return false;
}
```

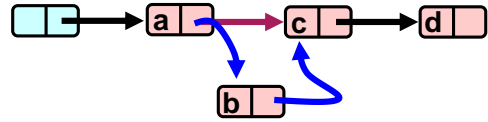
```
boolean contains(S, x) {
    int *curr = S.head;
    while(curr.key < x)
        curr = curr.next;
    if(curr.key == x)
        return true;
    return false;
}
```

```
typedef struct {
    int key;
    node *next;
} node;
```

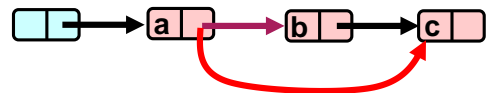
37

Sequential Operations

add ()



remove ()



38

Concurrent Sets

- What can happen if multiple threads call set operations at the "same time"?
 - Operations can conflict!
- Which operations conflict?
 - (add, remove), (add, add), (remove, remove), (remove, contains) will conflict
 - (add, contains) may miss update (which is fine)
 - (contains, contains) does not conflict
- How can we fix it?

39

Coarse-grained Locking

```
boolean add(S, x) {
    lock(S);
    node *pred = S.head;
    node *curr = pred.next;
    while(curr.key < x) {
        pred = curr;
        curr = pred.next;
    }
    if(curr.key == x)
        unlock(S);
        return false;
    else {
        node node = malloc();
        node.key = x;
        node.next = curr;
        pred.next = node;
    }
    unlock(S);
    return true;
}
```

```
boolean remove(S, x) {
    lock(S);
    node *pred = S.head;
    node *curr = pred.next;
    while(curr.key < x) {
        pred = curr;
        curr = pred.next;
    }
    if(curr.key == x) {
        pred.next = curr.next;
        free(curr);
        return true;
    }
    unlock(S);
    return false;
}
```

```
boolean contains(S, x) {
    lock(S);
    int *curr = S.head;
    while(curr.key < x)
        curr = curr.next;
    if(curr.key == x) {
        unlock(S);
        return true;
    }
    unlock(S);
    return false;
}
```

40

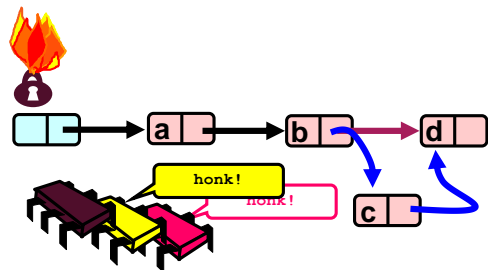
Coarse-grained Locking

- Correctness proof?
 - Assume sequential version is correct
 - Alternative: define set of invariants and proof that initial condition as well as all transformations adhere (pre- and post conditions)
 - Proof that all accesses to shared data are in CRs
 - This may prevent some optimizations
- Is the algorithm deadlock-free? Why?
 - Locks are acquired in the same order (only one lock)
- Is the algorithm starvation-free and/or fair? Why?
 - It depends on the properties of the used locks!

41

Coarse-grained Locking

- Is the algorithm performing well with many concurrent threads accessing it?



Simple but hotspot + bottleneck

42

Coarse-grained Locking

- Is the algorithm performing well with many concurrent threads accessing it?
 - No, access to the whole list is serialized
- BUT: it's easy to implement and proof correct
 - Those benefits should **never** be underestimated
 - May be just good enough
 - *"We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil. Yet we should not pass up our opportunities in that critical 3%. A good programmer will not be lulled into complacency by such reasoning, he will be wise to look carefully at the critical code; but only after that code has been identified"* — Donald Knuth (in *Structured Programming with Goto Statements*)

43

How to Improve?

- Will present some "tricks"
 - Apply to the list example
 - But often generalize to other algorithms
 - Remember the trick, not the example!
- See them as "concurrent programming patterns" (not literally)
 - Good toolbox for development of concurrent programs
 - They become successively more complex

44

Tricks Overview

1. Fine-grained locking
 - Split object into "lockable components"
 - Guarantee mutual exclusion for conflicting accesses to same component
2. Reader/writer locking
3. Optimistic synchronization
4. Lazy locking
5. Lock-free

45

Tricks Overview

1. Fine-grained locking
2. Reader/writer locking
 - Multiple readers hold lock (traversal)
 - contains() only needs read lock
 - Locks may be upgraded during operation
 - Must ensure starvation-freedom for writer locks!*
3. Optimistic synchronization
4. Lazy locking
5. Lock-free

46

Tricks Overview

1. Fine-grained locking
2. Reader/writer locking
3. Optimistic synchronization
 - Traverse without locking
 - Need to make sure that this is correct!*
 - Acquire lock if update necessary
 - May need re-start from beginning, tricky*
4. Lazy locking
5. Lock-free

47

Tricks Overview

1. Fine-grained locking
2. Reader/writer locking
3. Optimistic synchronization
4. Lazy locking
 - Postpone hard work to idle periods
 - Mark node deleted
 - Delete it physically later*
5. Lock-free

48

Tricks Overview

1. Fine-grained locking
2. Reader/writer locking
3. Optimistic synchronization
4. Lazy locking
5. Lock-free
 - Completely avoid locks
 - Enables wait-freedom
 - Will need atomics (see later why!)
 - Often very complex, sometimes higher overhead

49

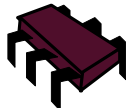
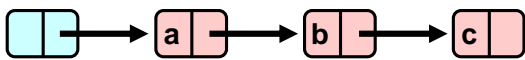
Trick 1: Fine-grained Locking

- Each element can be locked
 - High memory overhead
 - Threads can traverse list concurrently like a pipeline
- Tricky to prove correctness
 - And deadlock-freedom
 - Two-phase locking (acquire, release) often helps
- Hand-over-hand (coupled locking)
 - Not safe to release x's lock before acquiring x.next's lock
will see why in a minute
 - Important to acquire locks in the same order

```
typedef struct {
    int key;
    node *next;
    lock_t lock;
} node;
```

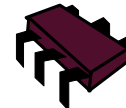
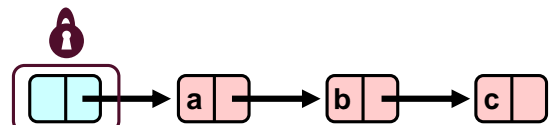
50

Hand-over-Hand (fine-grained) locking



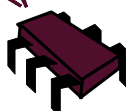
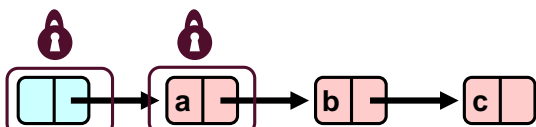
51

Hand-over-Hand (fine-grained) locking



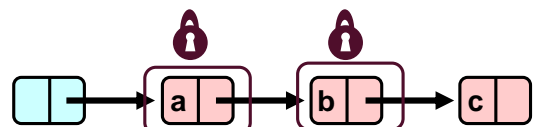
52

Hand-over-Hand (fine-grained) locking



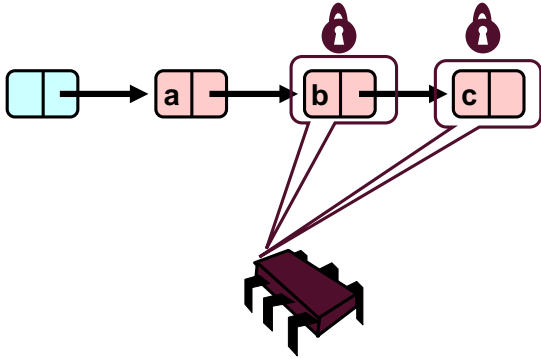
53

Hand-over-Hand (fine-grained) locking



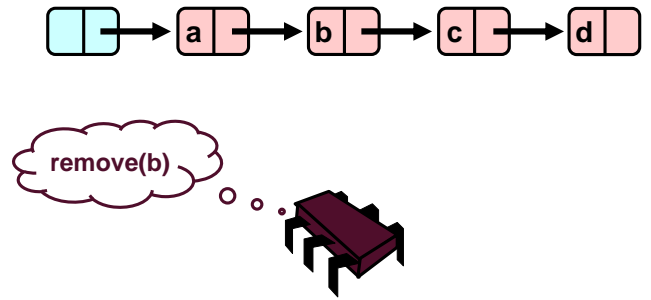
54

Hand-over-Hand (fine-grained) locking



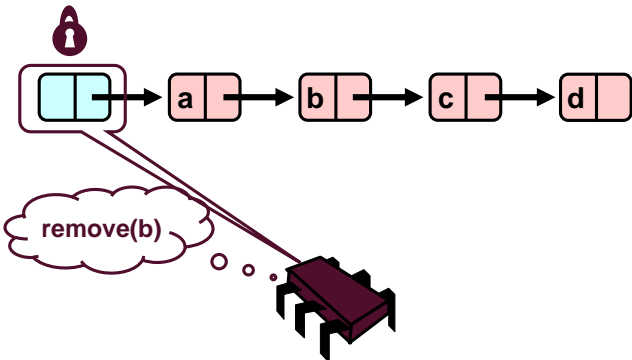
55

Removing a Node



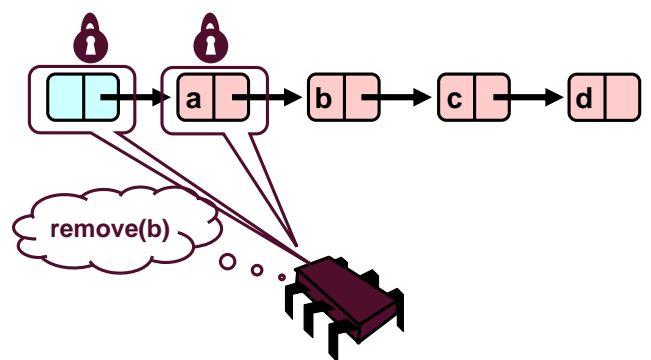
56

Removing a Node



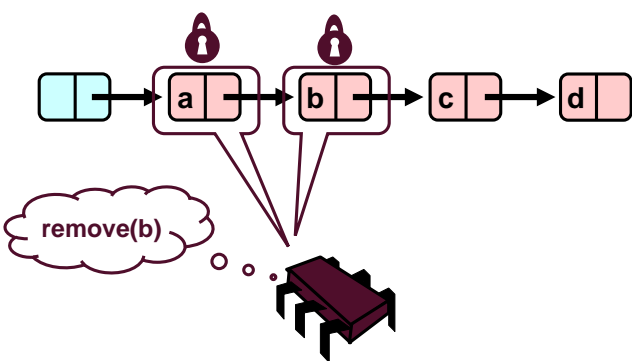
57

Removing a Node



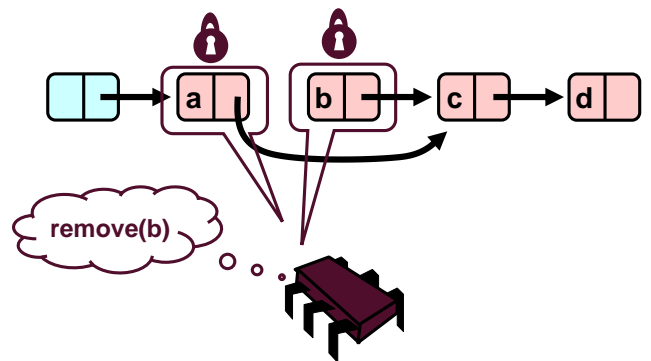
58

Removing a Node



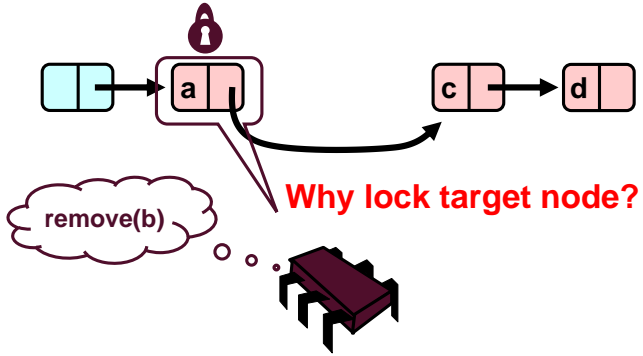
59

Removing a Node



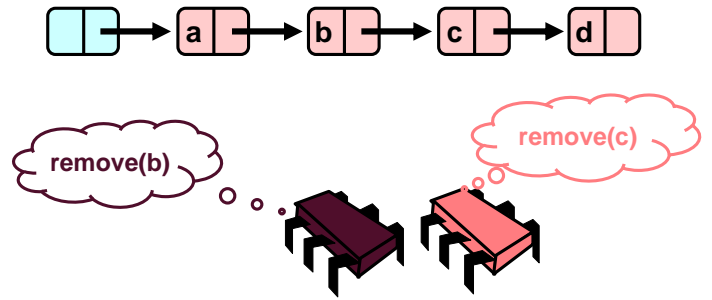
60

Removing a Node



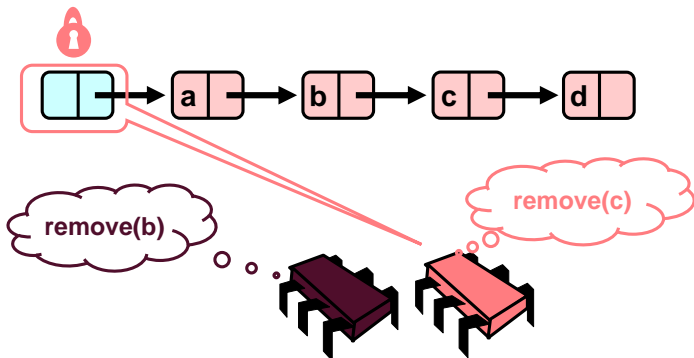
61

Concurrent Removes



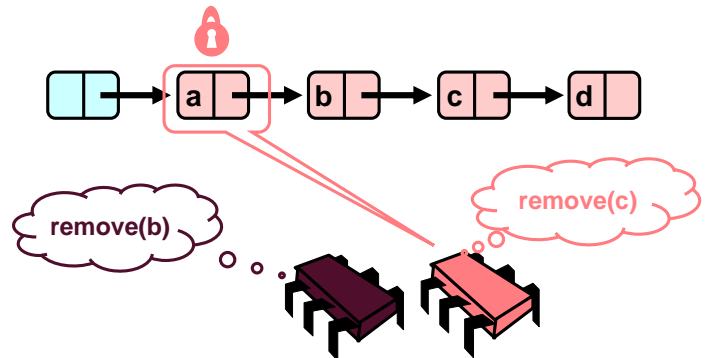
62

Concurrent Removes



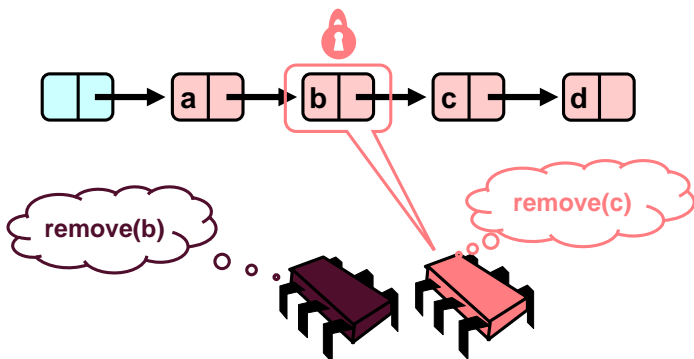
63

Concurrent Removes



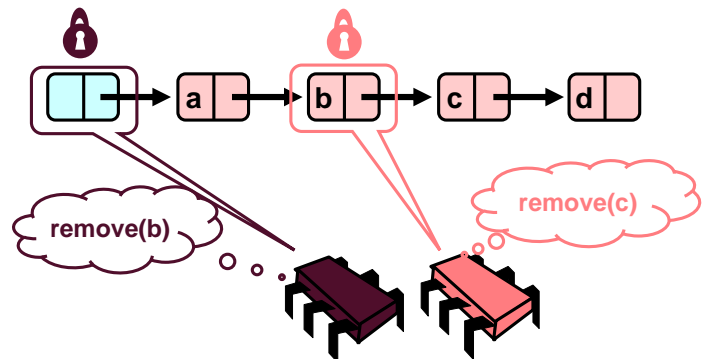
64

Concurrent Removes



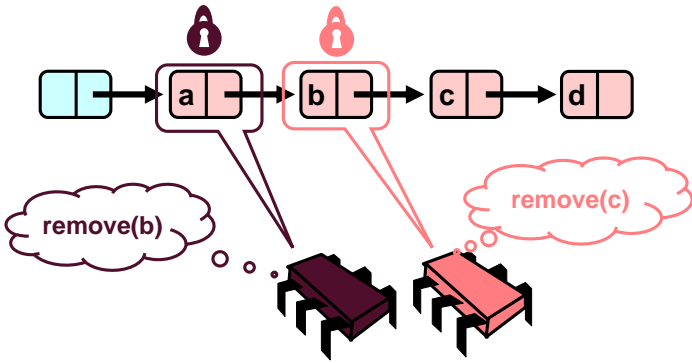
65

Concurrent Removes



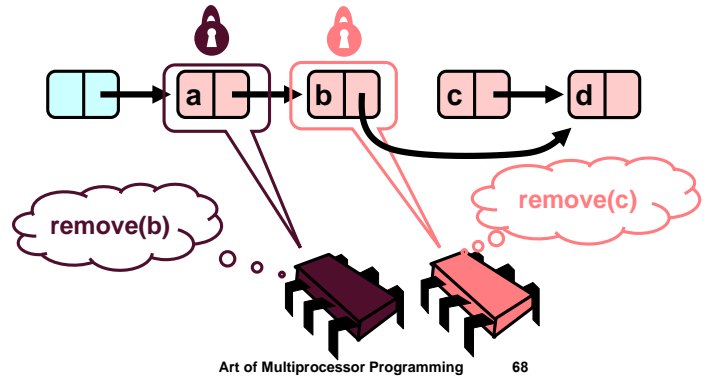
66

Concurrent Removes



67

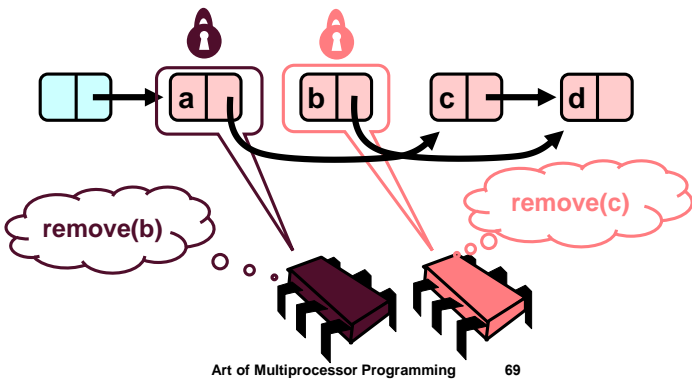
Concurrent Removes



Art of Multiprocessor Programming 68

68

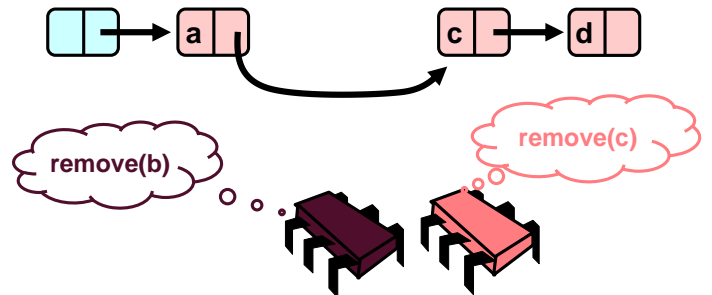
Concurrent Removes



Art of Multiprocessor Programming 69

69

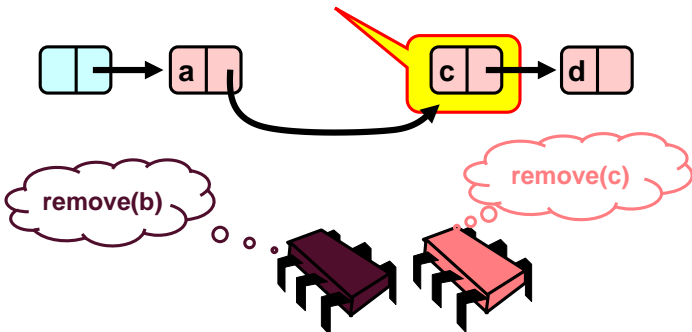
Uh, Oh



70

Uh, Oh

Bad news, c not removed



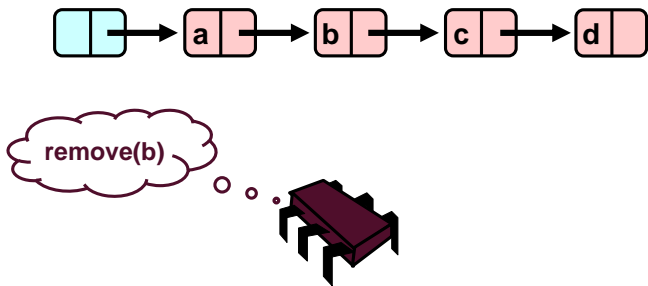
71

Insight

- If a node x is locked
 - Successor of x cannot be deleted!
- Thus, safe locking is
 - Lock node to be deleted
 - And its predecessor!
 - → hand-over-hand locking

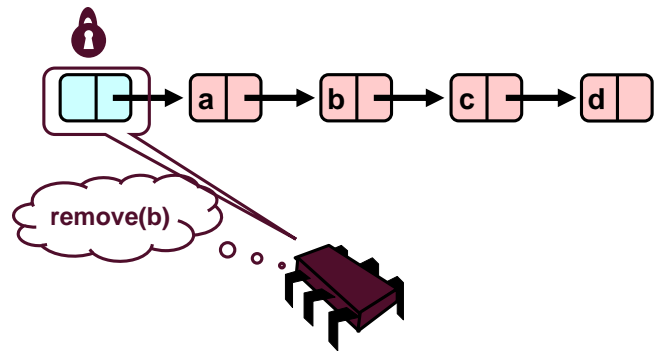
72

Hand-Over-Hand Again



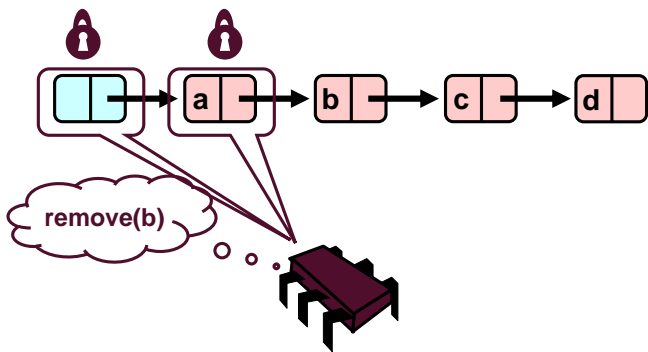
73

Hand-Over-Hand Again



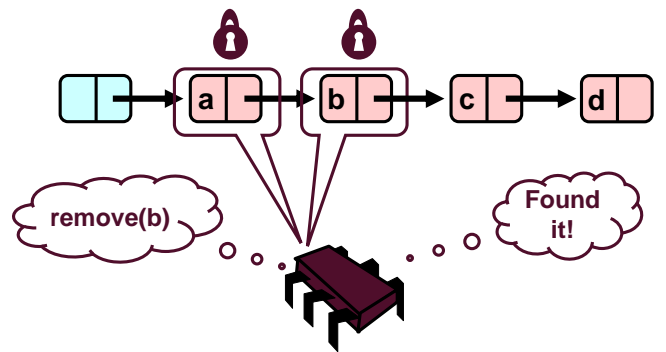
74

Hand-Over-Hand Again



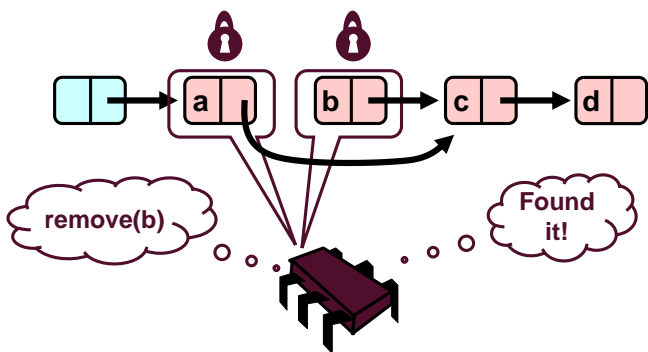
75

Hand-Over-Hand Again



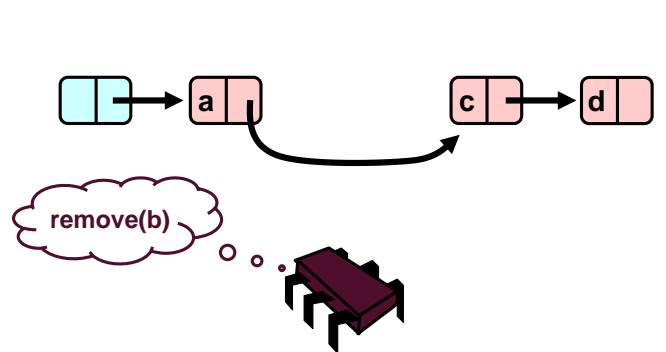
76

Hand-Over-Hand Again



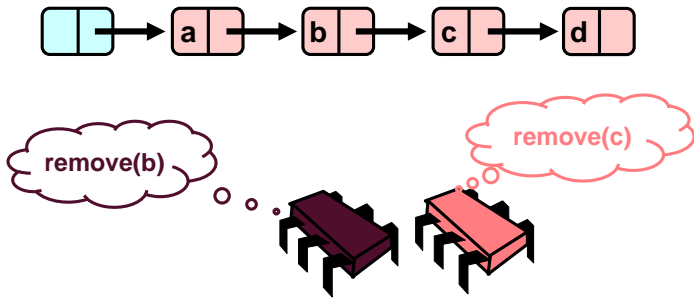
77

Hand-Over-Hand Again



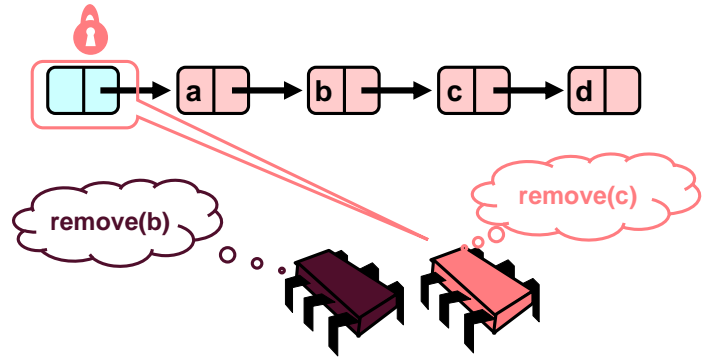
78

Removing a Node



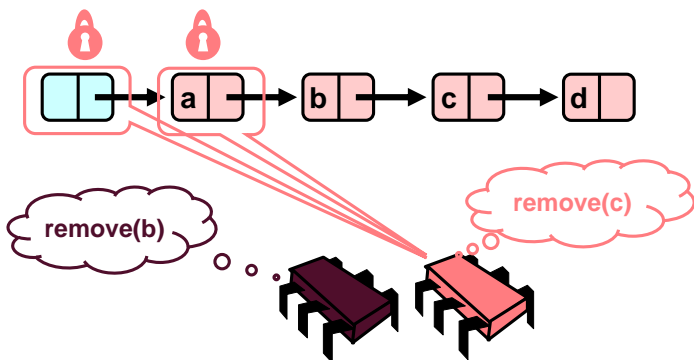
79

Removing a Node



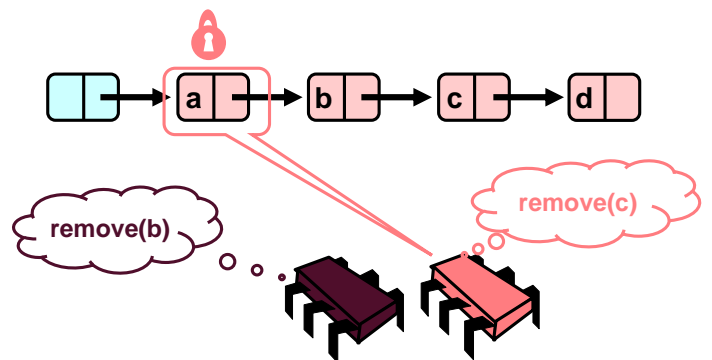
80

Removing a Node



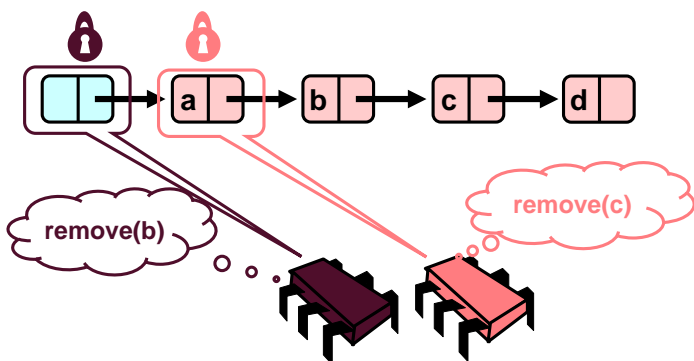
81

Removing a Node



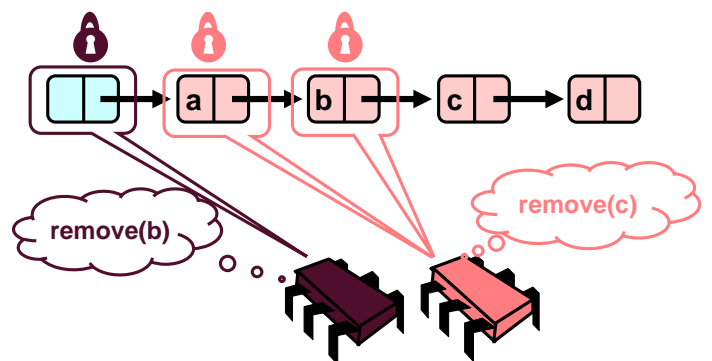
82

Removing a Node



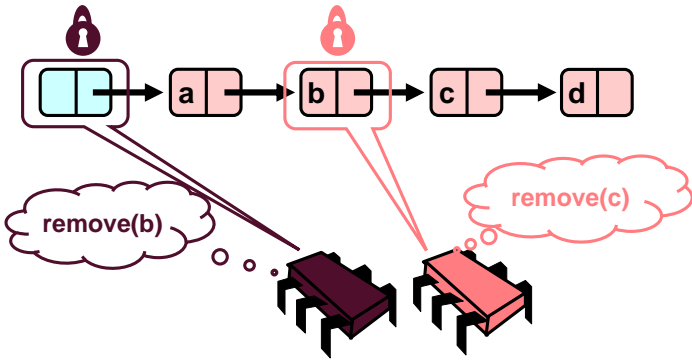
83

Removing a Node



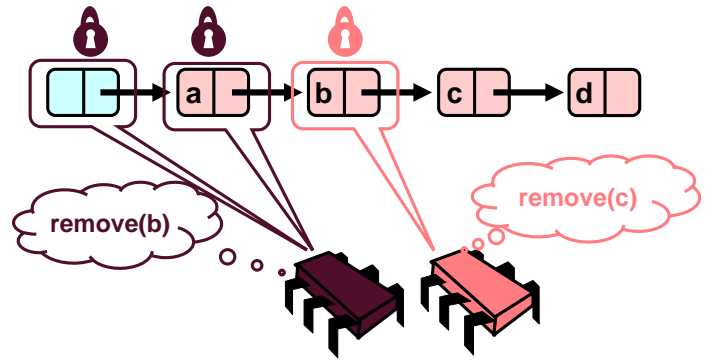
84

Removing a Node



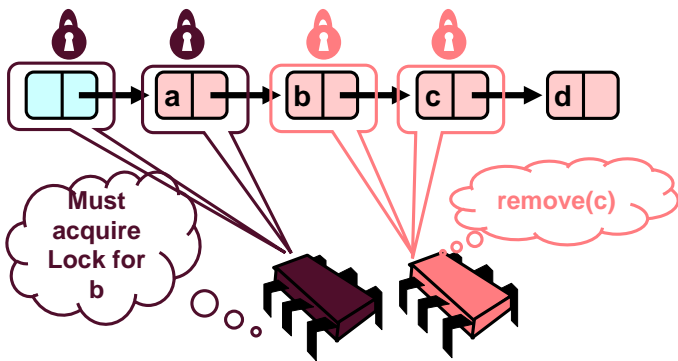
85

Removing a Node



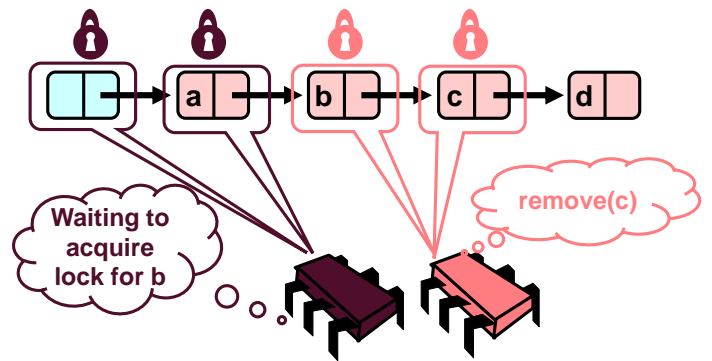
86

Removing a Node



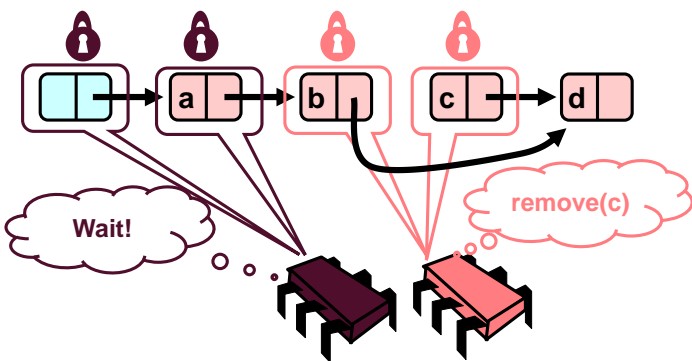
87

Removing a Node



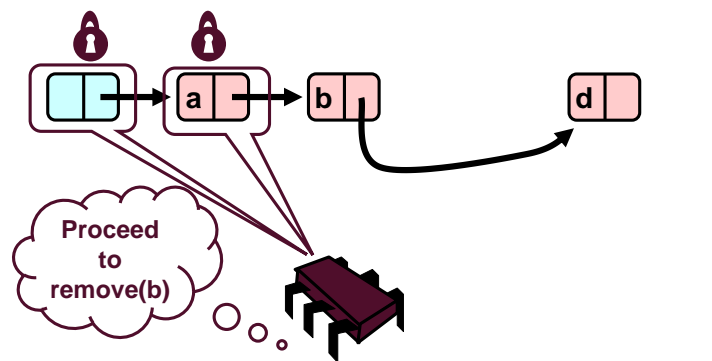
88

Removing a Node



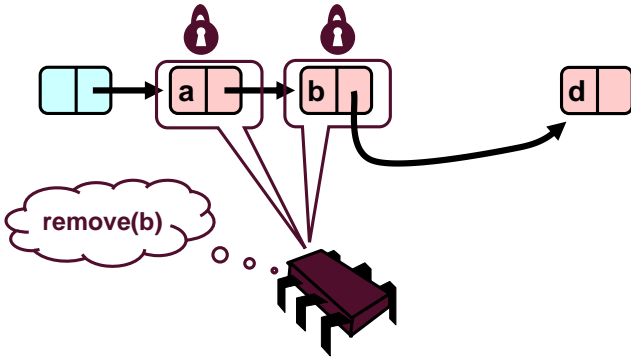
89

Removing a Node



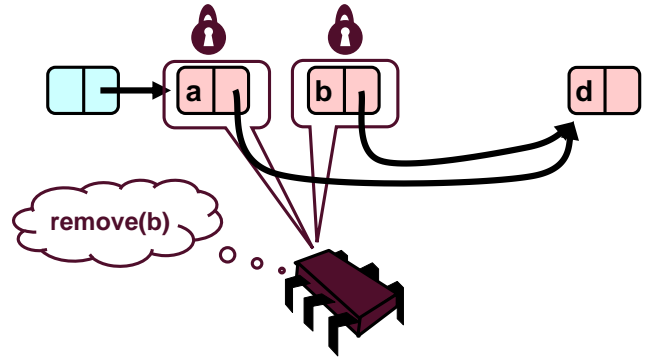
90

Removing a Node



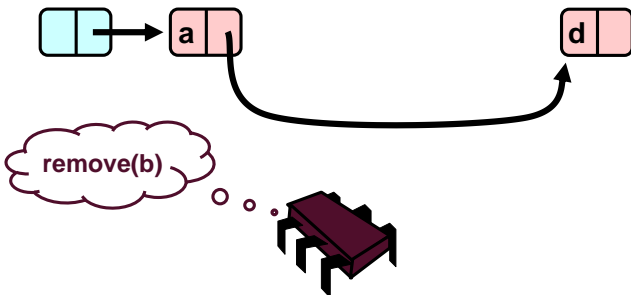
91

Removing a Node



92

Removing a Node



93

What are the Issues?

- **We have fine-grained locking, will there be contention?**
 - Yes, the list can only be traversed sequentially, a remove of the 3rd item will block all other threads!
 - This is essentially still serialized if the list is short (since threads can only pipeline on list elements)
- **Other problems, ignoring contention?**
 - Must acquire $O(|S|)$ locks

94

Trick 2: Reader/Writer Locking

- **Same hand-over-hand locking**
 - Traversal uses reader locks
 - Once add finds position or remove finds target node, upgrade **both** locks to writer locks
 - Need to guarantee deadlock and starvation freedom!
- **Allows truly concurrent traversals**
 - Still blocks behind writing threads
 - Still $O(|S|)$ lock/unlock operations

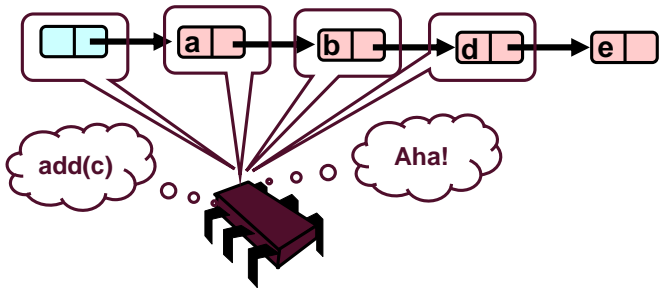
95

Trick 3: Optimistic synchronization

- **Similar to reader/writer locking but traverse list without locks**
 - Dangerous! Requires additional checks.
- **Harder to proof correct**

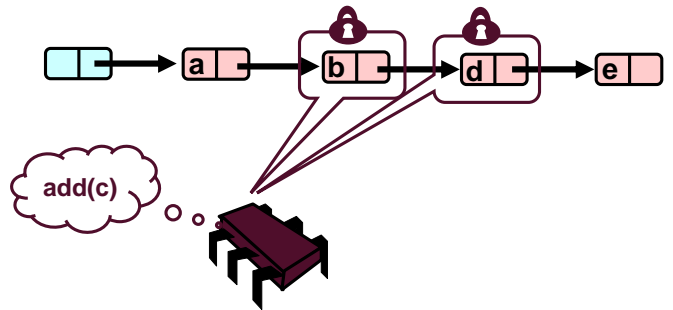
96

Optimistic: Traverse without Locking



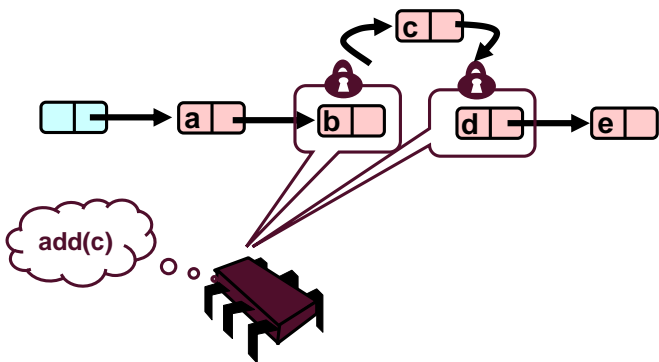
97

Optimistic: Lock and Load



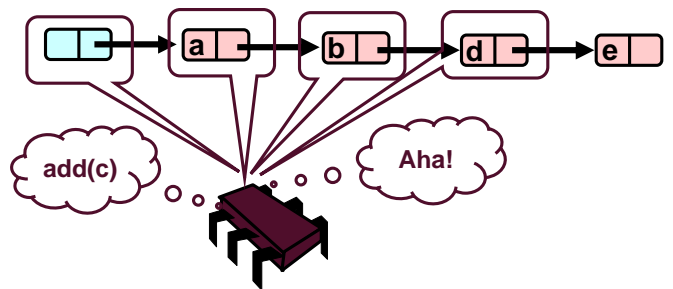
98

Optimistic: Lock and Load



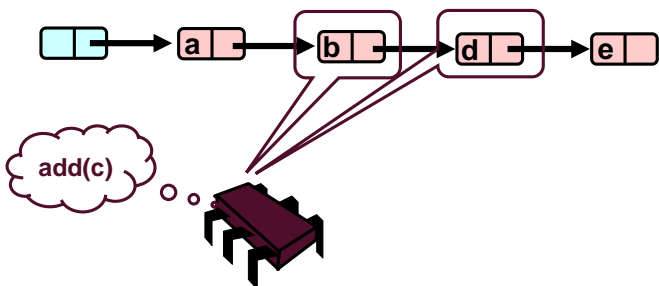
99

What could go wrong?



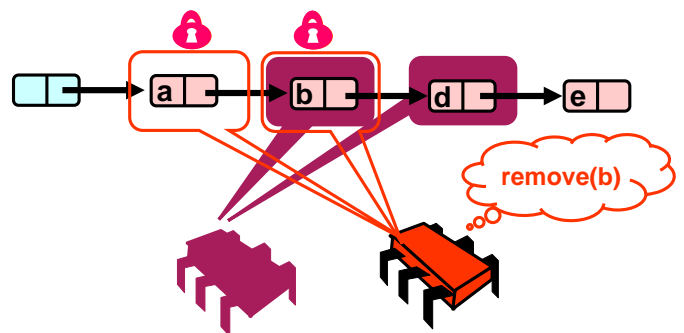
100

What could go wrong?



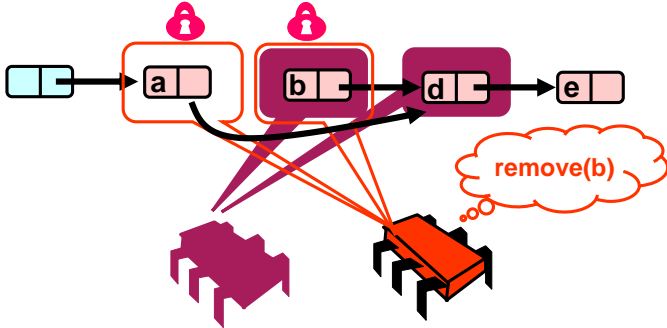
101

What could go wrong?



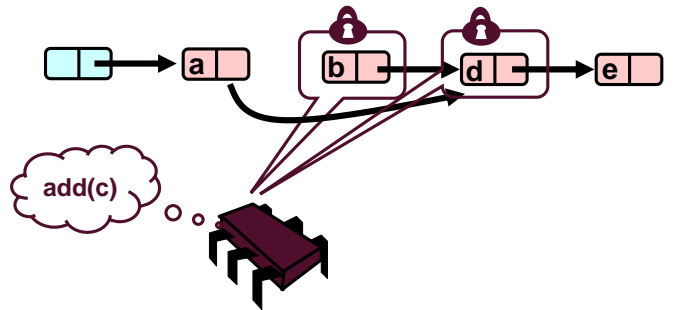
102

What could go wrong?



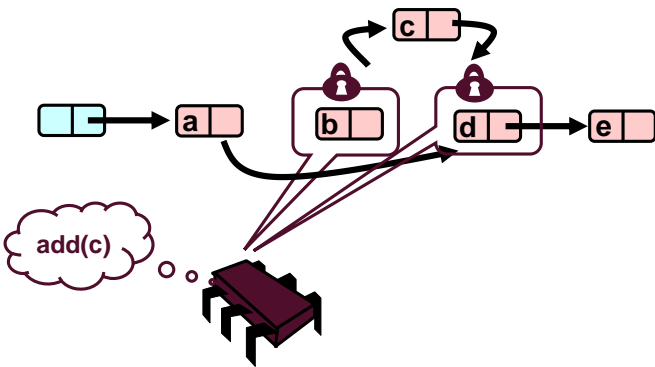
103

What could go wrong?



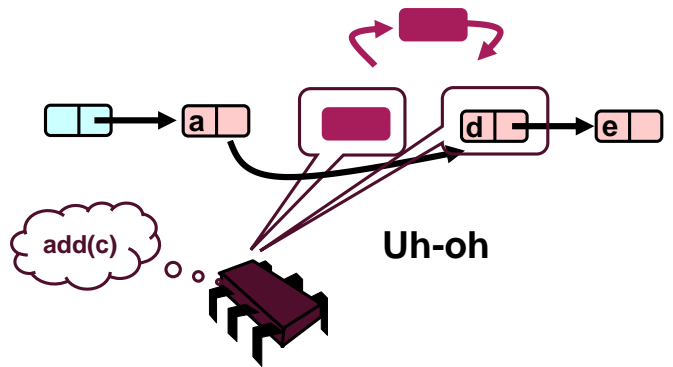
104

What could go wrong?



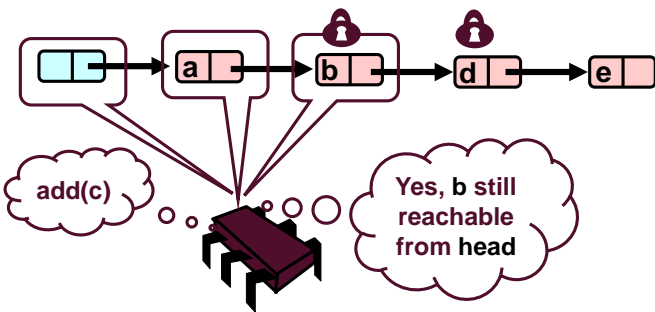
105

What could go wrong?



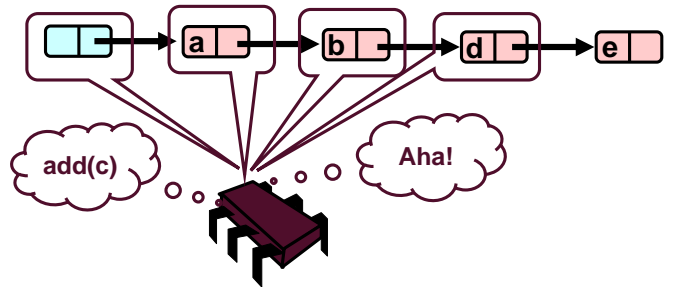
106

Validate – Part 1



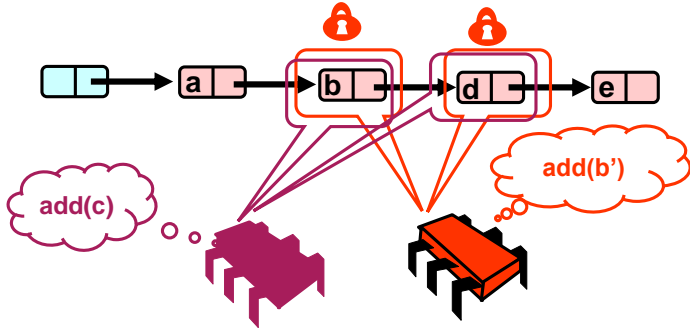
107

What Else Could Go Wrong?



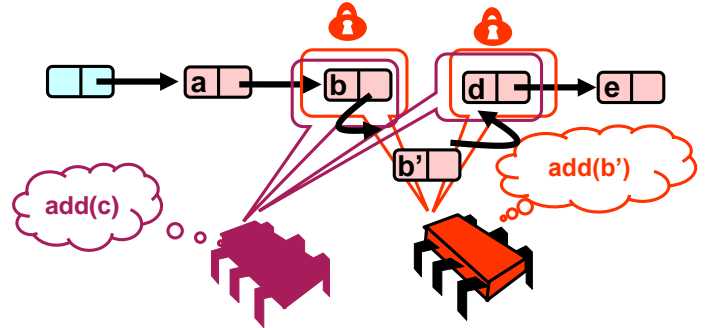
108

What Else Could Go Wrong?



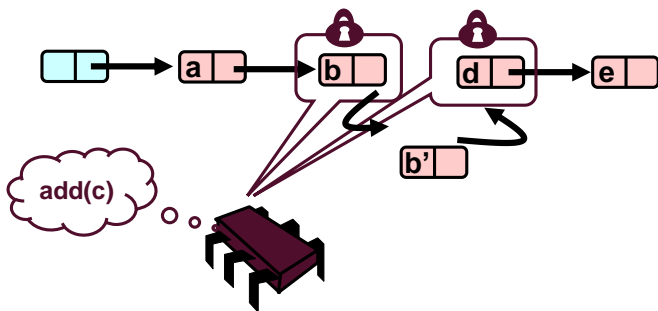
109

What Else Could Go Wrong?



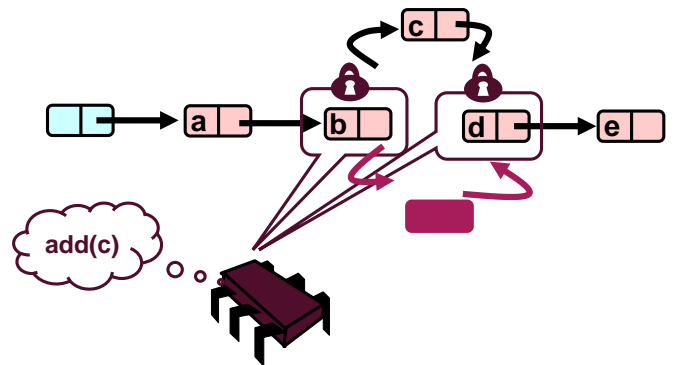
110

What Else Could Go Wrong?



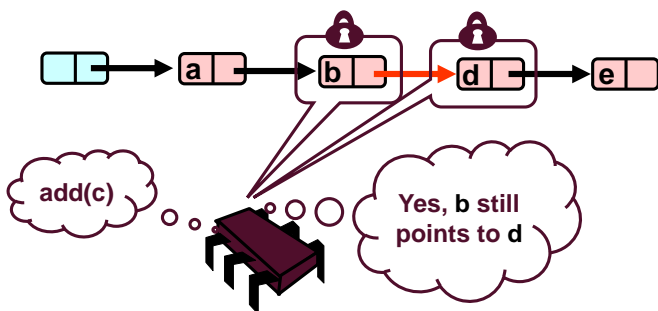
111

What Else Could Go Wrong?



112

Validate Part 2 (while holding locks)



113

Optimistic synchronization

- **One MUST validate AFTER locking**
 1. Check if the path how we got there is still valid!
 2. Check if locked nodes are still connected
 - If any of those checks fail?
 - Start over from the beginning (hopefully rare)*
- **Not starvation-free**
 - A thread may need to abort forever if nodes are added/removed
 - Should be rare in practice!
- **Other disadvantages?**
 - All operations requires two traversals of the list!
 - Even `contains()` needs to check if node is still in the list!

114

Trick 4: Lazy synchronization

- We really want one list traversal
- Also, `contains()` should be wait-free
 - Is probably the most-used operation
- Lazy locking is similar to optimistic
 - Key insight: removing is problematic
 - Perform it "lazily"
- Add a new "valid" field
 - Indicates if node is still in the set
 - Can remove it without changing list structure!
 - Scan once, `contains()` never locks!

```
typedef struct {
    int key;
    node *next;
    lock_t lock;
    boolean valid;
} node;
```

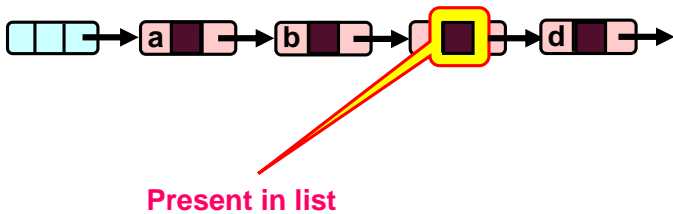
115

Lazy Removal



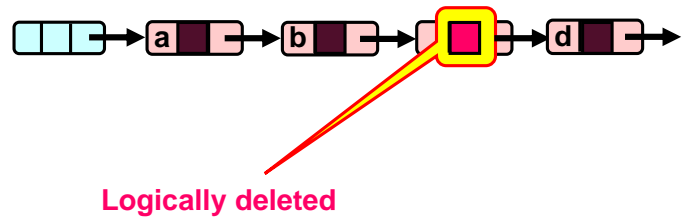
116

Lazy Removal



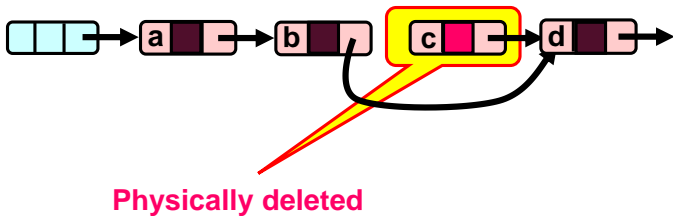
117

Lazy Removal



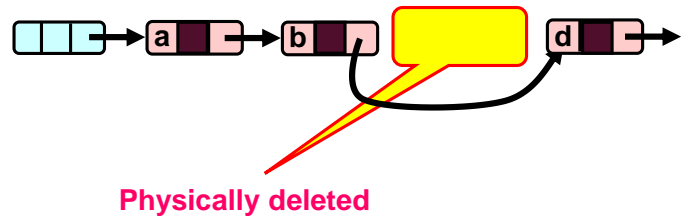
118

Lazy Removal



119

Lazy Removal



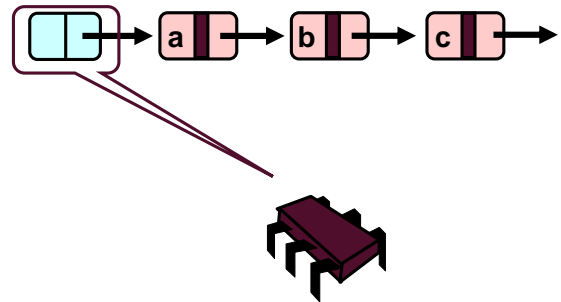
120

How does it work?

- **Eliminates need to re-scan list for reachability**
 - Maintains invariant that every **unmarked** node is reachable!
- **Contains can now simply traverse the list**
 - Just check marks, not reachability, no locks
- **Remove/Add**
 - Scan through locked and marked nodes
 - Removing does not delay others
 - Must only lock when list structure is updated
 - Check if neither pred nor curr are marked, pred.next == curr*

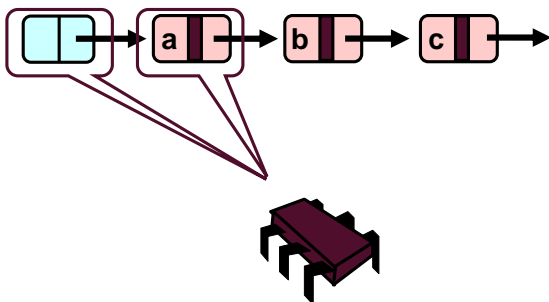
121

Business as Usual



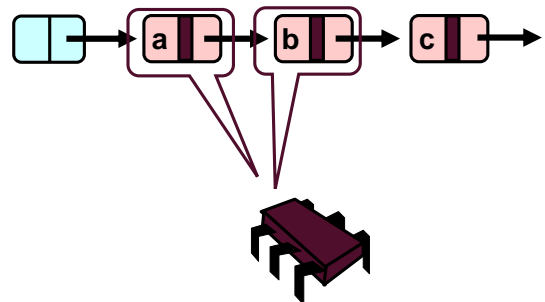
122

Business as Usual



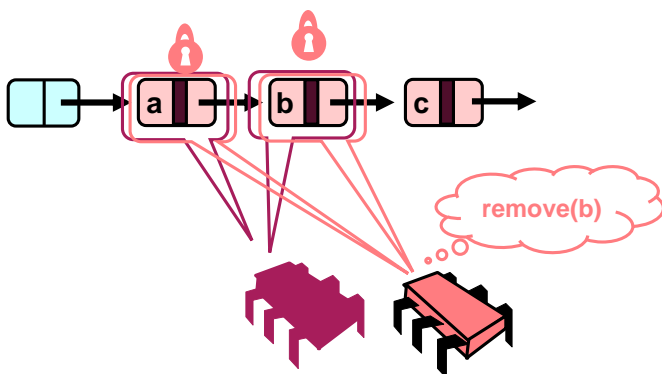
123

Business as Usual



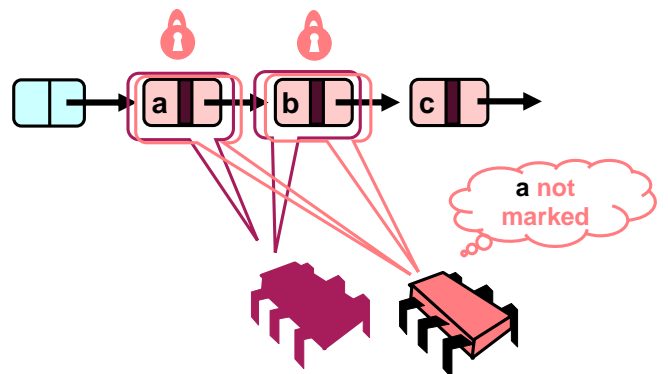
124

Business as Usual



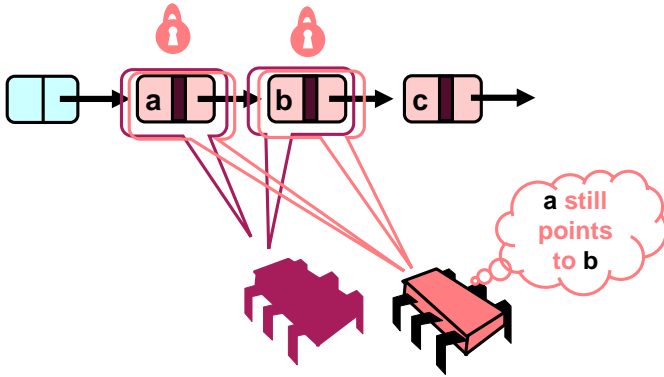
125

Business as Usual



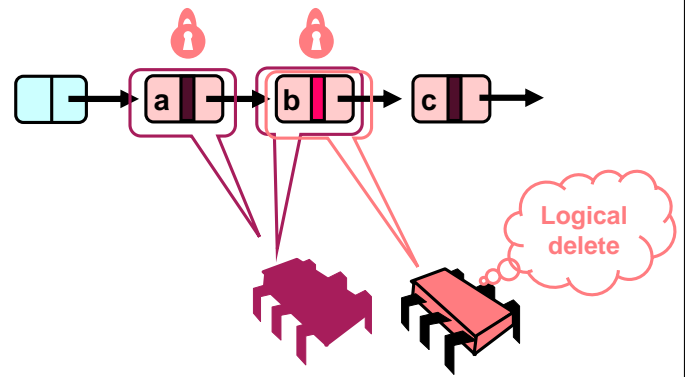
126

Business as Usual



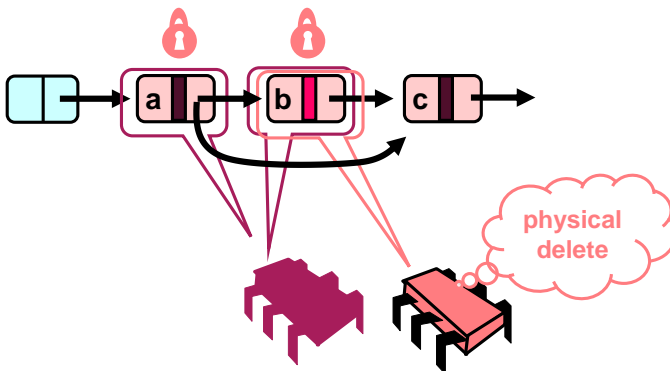
127

Business as Usual



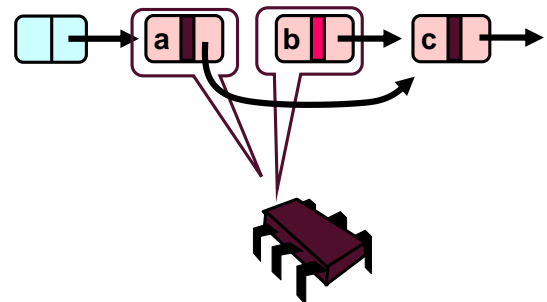
128

Business as Usual



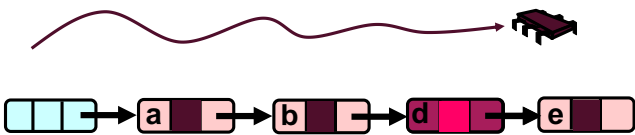
129

Business as Usual



130

Summary: Wait-free Contains



Use Mark bit + list ordering

1. Not marked → in the set
2. Marked or missing → not in the set

Lazy add() and remove() + Wait-free contains()

131

Problems with Locks

- What are the fundamental problems with locks?
- **Blocking**
 - Threads wait, fault tolerance
 - Especially when things like page faults occur in CR
- **Overheads**
 - Even when not contended
 - Also memory/state overhead
- **Synchronization is tricky**
 - Deadlock, other effects are hard to debug
- **Not easily composable**

132

Lock-free Methods

- **No matter what:**
 - Guarantee minimal progress
I.e., some thread will advance
 - Threads may halt at bad times (no CRs! No exclusion!)
I.e., cannot use locks!
 - Needs other forms of synchronization
E.g., atomics (discussed before for the implementation of locks)
Techniques are astonishingly similar to guaranteeing mutual exclusion

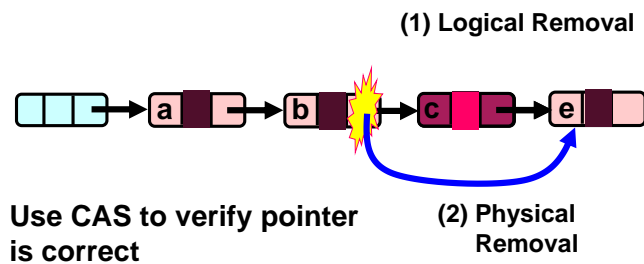
133

Trick 5: No Locking

- **Make list lock-free**
- **Logical succession**
 - We have wait-free contains
 - Make add() and remove() lock-free!
Keep logical vs. physical removal
- **Simple idea:**
 - Use CAS to verify that pointer is correct before moving it

134

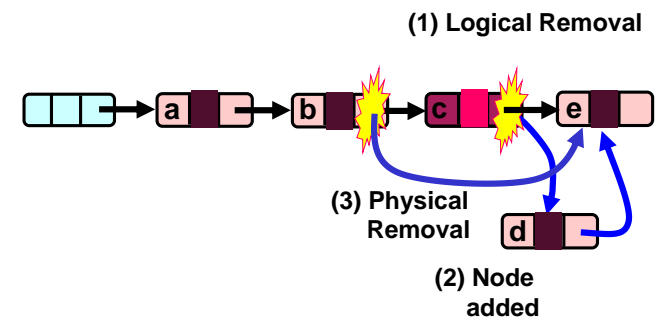
Lock-free Lists



Not enough! Why?

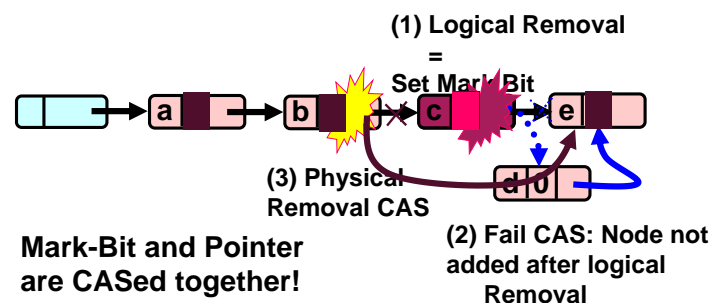
135

Problem...



136

The Solution: Combine Mark and Pointer



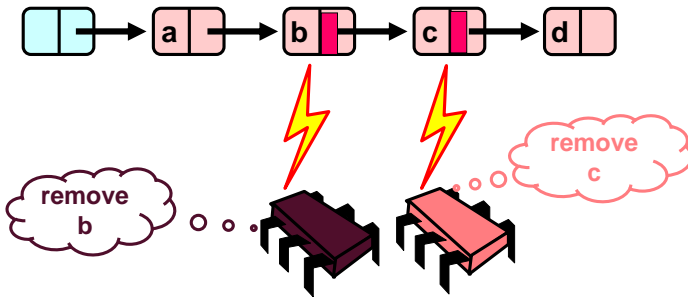
137

Practical Solution(s)

- **Option 1:**
 - Introduce "atomic markable reference" type
 - "Steal" a bit from a pointer
 - Rather complex and OS specific ☹
- **Option 2:**
 - Use Double CAS (or CAS2) ☹
CAS of two noncontiguous locations
 - Well, not many machines support it ☹
Any still alive?
- **Option 3:**
 - Our favorite ISA (x86) offers double-width CAS
Contiguous, e.g., lock cmpxchg16b (on 64 bit systems)
- **Option 4:**
 - TM!
E.g., Intel's TSX (essentially a cmpxchg64b (operates on a cache line))

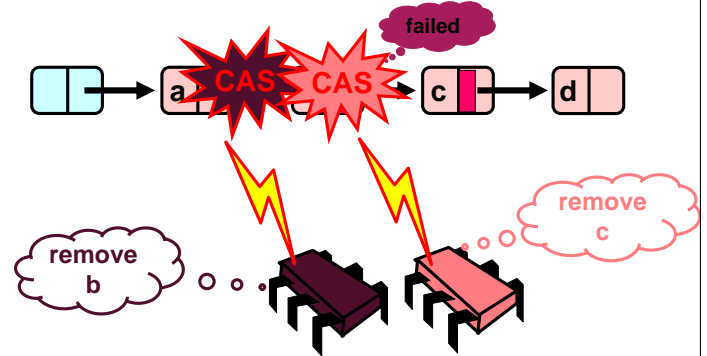
138

Removing a Node



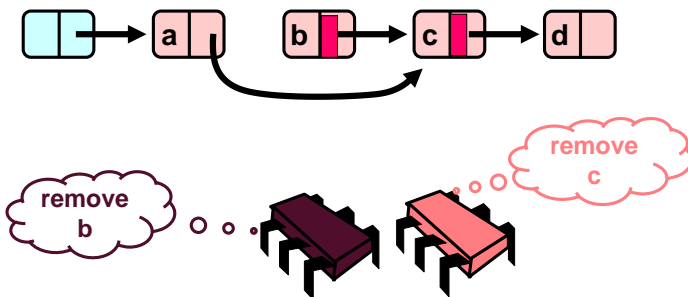
139

Removing a Node



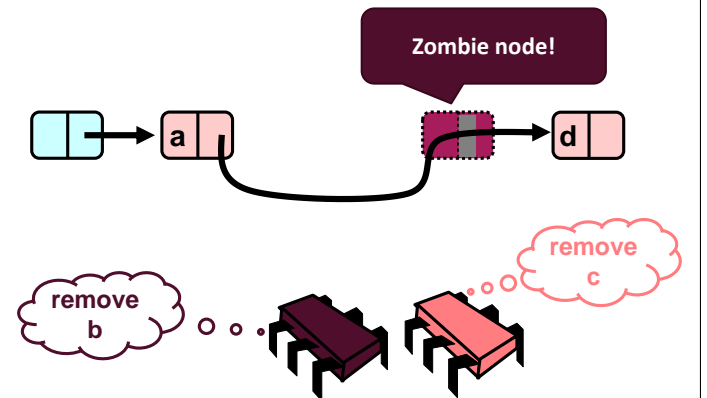
140

Removing a Node



141

Uh oh – node marked but not removed!



142

Dealing With Zombie Nodes

- **Add() and remove() "help to clean up"**
 - Physically remove any marked nodes on their path
 - I.e., if curr is marked: CAS (pred.next, mark) to (curr.next, false) and remove curr
 - *If CAS fails, restart from beginning!*
- "Helping" is often needed in wait-free algs
- This fixes all the issues and makes the algorithm correct!

143

Comments

- Atomically updating two variables (CAS2 etc.) has a non-trivial cost
- If CAS fails, routine needs to re-traverse list
 - Necessary cleanup may lead to unnecessary contention at marked nodes
- More complex data structures and correctness proofs than for locked versions
 - But guarantees progress, fault-tolerant and maybe even faster (that really depends)

144

More Comments

- **Correctness proof techniques**
 - Establish invariants for initial state and transformations
E.g., head and tail are never removed, every node in the set has to be reachable from head, ...
 - Proofs are similar to those we discussed for locks
Very much the same techniques (just trickier)
Using sequential consistency (or consistency model of your choice 😊)
Lock-free gets somewhat tricky
- **Source-codes can be found in Chapter 9 of “The Art of Multiprocessor Programming”**