CS 636: Concurrent Data Structures

Swarnendu Biswas

Department of Computer Science and Engineering, Indian Institute of Technology Kanpur

Sem 2024-25-II



Challenges with Concurrent Programming



Challenges with Concurrent Programming



Designing a Concurrent Set Data Structure

Designing a Concurrent Set

```
public interface Set<T> {
    boolean add(T x);
    boolean remove(T x);
    boolean contains(T x);
}
```

add(x)

adds x to the set and returns true if and only if x was not already present

remove(x)

removes x from the set and returns true if and only if x was present

contains(x)

returns true if and only if x is present in the set

There are significantly more calls to contains() than add() and remove()

Designing a Concurrent Set Using Linked Lists





Invariants

- Field key is data's hash code to help with efficient search
- Nodes are sorted based on the key value
- Assume that all hash codes are unique
- Sentinel nodes are immutable, and tail is reachable from head
- Removed nodes continue to represent valid memory locations

A Thread-Unsafe Set Data Structure

```
public class UnsafeList<T> {
    private Node head;
    public UnsafeList() {
        head = new Node(Integer.MIN_VALUE);
        head.next = new Node(Integer.
            MAX_VALUE);
    }
```

```
public boolean add(T x) {
  int key = x.hashcode();
  Node pred = head;
  Node curr = pred.next:
  while (curr.key < key) {</pre>
    pred = curr: curr = curr.next;
  if (kev == curr.kev) {
    return false;
  } else {
    Node node = new Node(x);
    node.next = curr:
    prev.next = node:
    return true;
```

2

0

11

12

13

14 15 16

A Thread-Unsafe Set Data Structure

```
public boolean remove(T x) {
     int kev = x.hashcode():
      Node pred = head:
     Node curr = pred.next;
     while (curr.kev < kev) {</pre>
        pred = curr;
        curr = curr.next;
8
     if (kev == curr.kev) {
9
        pred.next = curr.next:
        return true:
11
      } else {
12
        return false:
14
15
```

```
public boolean contains(T x) {
     int kev = x.hashcode():
     Node pred = head:
     Node curr = pred.next;
     while (curr.key < kev) {</pre>
       pred = curr;
       curr = curr.next;
     if (kev == curr.kev) {
       return true:
     } else {
       return false:
15
```

3

5

6

7

8

0

10

11

12

13

14

A Thread-Unsafe Set Data Structure



Thread-Unsafe Set: Incorrect remove()



Concurrent Set with Coarse-Grained Synchronization

```
public class CoarseList<T> {
                                                 20
  private Node head:
  private Lock lock = new ReentrantLock();
                                                 23
  public CoarseList() {
                                                 2/1
    head = new Node(Integer.MIN VALUE);
                                                 25
    head.next = new Node(Integer.MAX VALUE);
                                                 26
                                                 28
  . . .
  public boolean add(T x) {
                                                 30
    Node pred, curr
                                                 31
    int key = x.hashcode();
                                                 32
    lock.lock():
                                                 33
                                                 34
                                                 35
                                                 36
                                                 37
                                                 38
```

```
try {
  pred = head:
  curr = pred.next;
  while (curr.key < key) {</pre>
    pred = curr:
    curr = curr.next;
  if (kev == curr.kev) {
    return false;
  } else {
    Node node = new Node(x):
    node.next = curr:
    prev.next = node;
    return true;
} finally {
  lock.unlock():
```

2

9

10

12

14

15

16

18

19

Concurrent Set with Coarse-Grained Synchronization

```
public boolean remove(T x) {
  Node pred, curr;
  int key = x.hashcode();
  lock.lock():
  try {
    pred = head:
    curr = pred.next:
    while (curr.key < key) {</pre>
      pred = curr;
      curr = curr.next:
    if (kev == curr.kev) {
      pred.next = curr.next;
      return true:
    } else {
      return false;
  } finally {
    lock.unlock():
```

```
public boolean contains(T x) {
  Node curr;
  int key = x.hashcode();
  boolean found = false:
  lock.lock():
  trv {
    curr = head.next;
    while (curr.key < key) {</pre>
      curr = curr.next;
    if (kev == curr.kev) {
      found = true:
  } finally {
    lock.unlock():
  return found:
```

34

35

36

37

38

39

40

43

44

45

46

48

49

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

60

70

71

72

73

74

Concurrent Set with Fine-Grained Synchronization



Possible interleaving	
Thread 1	Thread 2
<pre>curr.lock.lock(); next = curr.next; curr.lock.unlock();</pre>	1 2 3
<pre>4 5 next.lock.lock();</pre>	4 // Remove next from list

Is Locking One Node Sufficient?



Concurrent Set with Fine-Grained Synchronization

```
public boolean add(T x) {
      int key = x.hashcode();
2
      head.lock();
      Node pred = head:
      try {
        Node curr = pred.next;
        curr.lock();
        trv {
8
          while (curr.key < key) {</pre>
0
            pred.unlock();
            pred = curr;
            curr = curr.next;
            curr.lock():
17.
15
```

```
if (kev == curr.kev) {
      return false:
    } else {
      Node node = new Node(x):
      node.next = curr;
      pred.next = node;
      return true;
  } finally {
    curr.unlock();
} finally {
  pred.unlock():
```

16

17

18

10

20

21

22

23

21.

25

26

27

28 29

30

Concurrent Set with Fine-Grained Synchronization

```
public boolean remove(T x) {
      int key = x.hashcode();
2
      head.lock();
     Node pred = null, curr = null;
     try {
        pred = head; curr = pred.next;
        curr.lock();
        trv {
8
          while (curr.key < key) {</pre>
0
            pred.unlock();
            pred = curr;
            curr = curr.next;
            curr.lock():
14
```

```
if (kev == curr.kev) {
      pred.next = curr.next;
      return true:
    } else {
      return false;
 } finally {
    curr.unlock():
} finally {
 pred.unlock();
```

15

16

17

18

19

20

21

22

23

24

25

26

27 28

Understanding Hand-over-Hand Locking

Hand-over-hand locking

The two locked nodes are always adjacent and guaranteed to be reachable from the head

Principles for efficient locking

- Acquire locks in a consistent order
- Do not hold too many locks
- Do not hold a lock for too long
- Only lock the written locations

Challenges With Fine-Grained Synchronization

Need to avoid deadlocks

- Deadlocks are always a problem with fine-grained locking
- For the Set data structure, each thread must acquire locks in some predetermined order

Are there other problems with the fine-grained Set design?

Challenges With Fine-Grained Synchronization

Need to avoid deadlocks

- Deadlocks are always a problem with fine-grained locking
- For the Set data structure, each thread must acquire locks in some predetermined order

Are there other problems with the fine-grained Set design?

- Potentially long sequence of lock acquire and release operations
- Prohibits concurrent accesses to disjoint parts of the data structure

Evaluating Concurrent Data Structures

Performance Metrics of Concurrent Data Structures

Speedup measures how effectively is an application utilizing resources

- Linear speedup is desirable
- Data structures whose speedup grow with resources is desirable

Amdahl's law says we need to reduce amount of serialized code

Reduce lock contention

Lock implementations with single memory location can introduce additional coherence and memory traffic due to unsuccessful acquires

Blocking or nonblocking implementations

Blocking Delay of any one thread can delay other threads

Nonblocking Delay of one thread cannot delay other threads

CS 636: Concurrent Data Structures

Reasoning about Correctness of Sequential Data Structures

Need to describe how an object's methods behave

- Possibilities include formal specification and API documentation
- Pre-condition describes the object's state before the method call
 - Operations on objects are not instantaneous. Each operation requires an invocation on that object, followed by a response.
 - ▶ Method call is the duration between an invocation event and a response event
- Post-condition describes the object's state and return value after the method call

Example

Suppose the state of a queue q is a sequence of items Q (i.e., precondition). Then, a call to $q \cdot enq(z)$ changes the state of the queue to $Q \cdot z$, where \bullet denotes concatenation.

Reasoning about Correctness of Concurrent Data Structures

Multiple threads can access a shared object, e.g., a node in our Set data structureSituation:Thread 1 is checking for contains(a)Thread 2 is executing remove(a)

Using pre- and post-conditions no longer work. How do you reason about the outcome?

We need ways to describe the correctness conditions for operations on a concurrent object

Reasoning about Correctness of Concurrent Data Structures

Correctness for interleaved operations on concurrent objects is determined by some notion of equivalence with sequential behavior

- Identify invariants and make sure they always hold
 - ▶ For example, an item is in the set if and only if it is reachable from head
- Correctness (or safety) property is linearizability
- Progress (or liveness) properties are starvation and deadlock-freedom

Definitions

Program order The order in which a single thread issues method calls is called its program order.

• Method calls by different threads are unrelated by program order.

Total method A method is total if it is defined for every object state, i.e., it does not need to wait for certain conditions to become true.

• A total method is used when the caller thread has something useful to do than wait for certain conditions to be met.

Partial method A partial method is not defined for every object state, it may have to block for certain conditions to hold.

• For example, a partial Queue::get() call that tries to remove an item from an empty queue blocks until an item is available to return.

Compositional A correctness property P is compositional if, whenever each object in the system satisfies P, the system as a whole satisfies P.

Correct Behavior Expected From a Concurrent Execution



Sequentially Consistent Execution

An execution is sequentially consistent (SC) if the result is the same as if the operations from all threads were executed in some sequential order, and the operations of each individual thread appear in program order.

Consider the following operations on a FIFO queue q, where x and y are objects. Thread 1 q.enq(x) q.deq(y) Thread 2 q.enq(y) q.deq(x)

There are two possible sequential orders that can justify the above execution

Swarnendu Biswas (IIT Kanpur)

SC can Violate Real-Time Order

Reordering method calls unrelated by program order is allowed in SC, and so can violate real-time order



SC is Not Composable

p and q are each sequentially consistent, but the execution as a whole is not



Linearizability



Swarnendu Biswas (IIT Kanpur)

Linearizability

Linearizability is popularly used to argue that a concurrent algorithm correctly implements a problem through a sequential specification

- Linearizability has two requirements
 - (i) Method calls should appear to happen one-at-a-time in sequential order
 - (ii) Each method call should appear to take effect instantaneously at some moment between its invocation and response
- Linearization point represents a single atomic step where the method call "takes effect"
 - For coarse-grained lock-based implementations, each method's critical section is its linearization point
 - ► For implementations that do not use locking, the linearization point is a single step where the effects of the method call become visible to others

Concurrent Operations on an Object

- Say you perform some operations on an object (e.g., a method call)
- A history is a sequence of invocations and responses on an object made by concurrent threads



Sequential History

- Sequential history is where all invocations and responses are instantaneous
 - ► Starts with an invocation, last invocation may not have a response
 - Method calls do not overlap



Sequential History

- Sequential history is where all invocations and responses are instantaneous
 - ► Starts with an invocation, last invocation may not have a response
 - Method calls do not overlap



Linearizable History

- Every concurrent history is equivalent to some sequential history
 - If one method call precedes another, then the earlier call must have taken effect before the later call
 - ▶ If two method calls overlap, we can order them in any way
- Consider a concurrent history (set of method calls) *H* and a valid sequential history *S*. The history *H* is linearizable if:
 - ► For every completed call in *H*, the call returns the same result as it would return if every operation in *H* would have been completed one after the other (i.e., in *S*)
 - ▶ If method call m_1 completes before method call m_2 in H, then m_1 precedes m_2 in S

Linearizability in Simpler Words

- Sequential history is correct according to the semantics of the object
- Invocations and responses can be reordered to form a sequential history
- If a response preceded an invocation in the original history, it must still precede it in the sequential reordering
Understanding Linearizability



Understanding Linearizability



Identifying Linearization Points

Linearization point represents a single atomic step where the method call "takes effect", and is between the function invocation and response

What are the linearization points for the methods add(), remove(), and contains() for the coarsely- and finely-synchronized Set?

Sequential Consistency vs Linearizability

Sequential Consistency

- Method calls appear to happen instantaneously in some sequential order
- A sequentially consistent history is not necessarily linearizable
- Nonblocking but not composable

Linearizability

- Method calls appear to happen instantaneously at some point between its invocation and response
- Every linearizable history is sequentially consistent
- Nonblocking and composable

Progress Guarantees

wait-free A method is wait-free if it guarantees that every call finishes in a finite number of steps

- lock-free A method is lock-free if it guarantees that some call always finishes in a finite number of steps
- obstruction-free A method is obstruction-free if it is guaranteed to finish its operations in the absence of contention

Designing a Concurrent Set Data Structure

How to Design a Concurrent Set?

Coarse-grained synchronization

Easy to get right, low concurrency, not scalable

Fine-grained synchronization

More concurrent and scalable than coarse-grained synchronization, difficult to get right

Optimistic synchronization

Avoid synchronization to search, good for low contention cases

Lazy synchronization

Defer expensive data structure manipulation operations

Nonblocking synchronization

Rely on atomic operations such as compareAndSet() for synchronization

Swarnendu Biswas (IIT Kanpur)

CS 636: Concurrent Data Structures

Optimistic Synchronization

Optimistic strategy

- Access data without acquiring a lock
- Lock only when required, and validate that the condition before locking is still valid
- If valid, then continue with access/mutation
- If invalid, restart by locking again

Optimistic strategy works well if conflicts are rare

```
public boolean add(T x) {
                                                 17
      int key = x.hashcode();
2
                                                 18
      while (true) {
                                                 19
         Node pred = head:
                                                 20
         Node curr = pred.next;
                                                 21
         while (curr.kev < kev) {</pre>
                                                 22
           pred = curr;
                                                 23
           curr = curr.next;
8
                                                 24
                                                 25
0
         pred.lock(); curr.lock();
                                                 26
                                                 27
                                                 28
                                                 29
                                                 30
17.
15
                                                 31
16
                                                 32
```

```
try {
  if (validate(pred, curr)) {
    if (curr.key == key) {
      return false:
    } else {
      Node node = new Node(x):
      node.next = curr;
      prev.next = node;
      return true;
} finally {
  curr.unlock(): pred.unlock():
```

Is Validation Necessary?

Thread 1 is executing remove(p)

Other threads execute remove(b-p)



How to Validate?

Double check that the optimistic result is still valid

```
// Check that prev is reachable from head and prev.next == curr
boolean validate(Node prev, Node curr) {
    Node node = head;
    while (node.key <= prev.key) {
        if (node == prev)
            return prev.next == curr;
        node = node.next;
    }
    return false;
}
```

```
public boolean remove(T x) {
     int key = x.hashcode();
2
     while (true) {
        Node pred = head:
        Node curr = pred.next;
        while (curr.kev < kev) {</pre>
          pred = curr;
          curr = curr.next;
8
0
        pred.lock(); curr.lock();
17.
```

```
try {
  if (validate(pred, curr)) {
    if (curr.kev == kev) {
      pred.next = curr.next;
      return true;
    } else {
      return false;
} finally {
  curr.unlock(); pred.unlock();
```

15

16

17

18

19

20

21

22 23

24



Are there problems with the optimistic design?

- Validation can be costly (e.g., need to traverse the list again)
- Needs lock operations for contains() which is the most frequent method

Lazy Synchronization

Delay mutation operations for a later time

- Add a mark or flag bit on each node to indicate logical deletion
- Invariant: every unmarked node is reachable from head

Guarantees

add() traverses the list, locks the predecessor, and inserts the node remove() marks the target node logically removing it, then redirects the predecessor's next link physically removing the target node contains() needs only one wait-free traversal (no locking is required)

Concurrent Set with Lazy Synchronization

```
} else {
   public boolean add(T x) {
                                              17
                                                               Node node = new Node(x);
     int key = x.hashcode();
2
                                              18
     while (true) {
                                                               node.next = curr;
3
                                              19
       Node pred = head;
                                                               pred.next = node;
                                              20
       Node curr = pred.next;
                                                               return true;
                                              21
       while (curr.kev < kev) {</pre>
          pred = curr;
                                              23
                                                        } finally {
         curr = curr.next:
8
                                              24
       }
                                                          curr.unlock(); }
                                              25
0
       pred.lock();
                                              26
       trv {
                                                      } finally {
11
                                              27
                                                        pred.unlock();
         curr.lock();
                                              28
          trv {
                                              29
            if (validate(pred, curr)) {
                                              30
17.
              if (curr.key == key) {
                                              31
15
                return false:
16
                                              32
```

How to Validate?

Check that both prev and curr are unmarked and prev.next == curr

boolean v	validate(Node	pre	v, Node curr) {			
return	<pre>!prev.marked</pre>	66	<pre>!curr.marked</pre>	ծծ	prev.next	==	curr;
}							

Concurrent Set with Lazy Synchronization

```
public boolean remove(T x) {
                                                             } else {
                                             17
      int key = x.hashcode();
                                                               // Logical deletion
2
                                             18
      while (true) {
                                                               curr.marked = true;
                                             19
                                                               // Physical deletion
        Node pred = head:
                                             20
        Node curr = pred.next;
                                                               pred.next = curr.next;
                                             21
        while (curr.kev < kev) {</pre>
                                                               return true;
                                             22
          pred = curr;
                                             23
          curr = curr.next;
8
                                             24
                                                        } finally {
                                             25
0
        pred.lock();
                                                           curr.unlock(); }
                                             26
        trv {
          curr.lock();
                                                      } finally {
                                             28
                                                        pred.unlock():
          trv {
                                             29
             if (validate(pred, curr)) {
                                             30
17.
               if (curr.key != key) {
                                             31
15
                 return false:
16
                                             32
```

Detecting Conflicts: Scenario 1

Thread 1 is executing remove(b)

Thread 2 is executing remove(a)



Detecting Conflicts: Scenario 2

Thread 1 is executing remove(b)

Thread 2 is executing add(p)



Concurrent Set with Lazy Synchronization

```
public boolean contains(T x) {
    int key = x.hashcode();
    Node curr = head;
    while (curr.key < key) {
        curr = curr.next;
     }
    return curr.key == key && !curr.marked;
   }
</pre>
```



Unsuccessful contains(): Scenario 1



CS 636: Concurrent Data Structures

Sem 2024-25-II

Unsuccessful contains(): Scenario 2



Thread 1 is traversing along the marked portion of the list $p \dots x$

Swarnendu Biswas (IIT Kanpur)

CS 636: Concurrent Data Structures

Unsuccessful contains(): Scenario 2

Thread 1 is executing contains(x)

Thread 2 is executing add(x)



Thread 1 is traversing along the marked portion of the list $p \dots x$

Nonblocking Synchronization

Why do we need nonblocking designs?

- Use of locks can lead to deadlocks, livelocks, and priority inversion
- Blocked threads do not do useful work, problematic for high-priority or real-time applications
- Getting the right degree of concurrency and correctness with locks is challenging

Idea: Use RMW instructions like CAS to update the next field

Eliminate locks altogether

Nonblocking Algorithms

- + Failure or suspension of a thread does not impact other threads
- + Guaranteed system-wide progress implies lock-freedom, while per-thread progress implies wait-freedom
 - ▶ Wait-freedom is the strongest nonblocking progress guarantee
 - Lock-freedom allows an individual thread to starve
 - All wait-free algorithms are lock-free
- Lock-free implies "locking up" the application in some way (e.g., deadlock and livelock)
 - ► Lock-free does **not only** imply absence of synchronization locks

Compare-and-Swap (CAS) Primitive

- Modern architectures provide many atomic read-modify-write (RMW) instructions for synchronization
 - For example, test-and-set, fetch-and-add, compare-and-swap, and load-linked/store-conditional
- Compare-and-Swap (CAS) compares the contents of a memory location with a given value and, only if they are the same, updates the contents of that memory location to a new given value

```
bool CAS(word* loc, word oldval, word newval) {
    atomic { // Code block will execute atomically
    res := (*loc == oldval);
    if (res)
      *loc := newval;
    return res;
    }
}
```

2

3

4

5

Compare-and-Swap (CAS) Primitive

- CAS is implemented as the compare-and-exchange (CMPXCHG) instruction in x86 architectures
 - On a multiprocessor, the LOCK prefix must be used
- CAS is a popular synchronization primitive for implementing both lock-based and nonblocking concurrent data structures

1	xor %ecx, %ecx /*ecx=0*/	1	<pre>void spinlock(lock* lk) {</pre>
2	inc $%ecx / *ecx = 1 * /$	2	// flag attribute is set when the
		2	// lock is acquired
3 KEIKT:	XUP $\%$ edX, $\%$ edX $/*$ edX $= 0*/$	3	// LOCK IS acquired
4	lock compxcng %ecx, block	4	while $(CAS(SIK->flag, 0, 1) == 1)$ {
5	jnz RETRY	5	// Keep spinning
6	ret	6	}
7		7	}

Nonblocking Synchronization with CAS on Only next Field



a is deleted but b is not added to the list

Swarnendu Biswas (IIT Kanpur)

CS 636: Concurrent Data Structures

Sem 2024-25-II

Nonblocking Synchronization with CAS on Only next Field

Thread 1 is executing remove(a)

Thread 2 is executing remove(b)



a is deleted but b is not deleted from the list

Swarnendu Biswas (IIT Kanpur)

CS 636: Concurrent Data Structures

Sem 2024-25-II

60 / 112

Disallow Updates to Deleted Nodes

- Cannot allow updates to a node once it has been logically or physically removed from the list
- Treat the next and marked fields as atomic
 - ▶ An attempt to update the next field when the marked field is true will fail

Java provides Class AtomicMarkableReference<T> in the java.util.concurrent.atomic package

```
public boolean compareAndSet(T expectedReference, T newReference, boolean
        expectedMark, boolean newMark);
public T get(boolean[] marked);
public T getReference();
public Boolean isMarked();
```

Designing a Nonblocking Set

- The next field is of type AtomicMarkableReference<Node>
- A thread logically removes a node by setting the marked bit in the next field
- As threads traverse the list, they clean up the list by physically removing marked nodes
 - Threads performing add() and remove() do NOT traverse marked nodes, they remove them before continuing



Challenge in Traversing Marked Nodes

Thread 1 is executing remove(b)

Thread 2 marks a



Challenge in Traversing Marked Nodes

Thread 1 is executing remove(b) Thread 2 marks a



Thread 1 does not delete the marked node a \implies Thread 1 cannot redirect a.next

Swarnendu Biswas (IIT Kanpur)

CS 636: Concurrent Data Structures

Sem 2024-25-II

63 / 112

Helper method public Window find(Node head, int key)

Traverses the list seeking to set pred to the node with the largest key less than key, and curr to the node with the least key greater than or equal to key

```
class Window {
   public Node pred, curr;
   Window(Node myPred, Node myCurr) {
      pred = myPred; curr = myCurr;
   }
}
```

Helper Method

```
public Window find(Node head, int key) {
      Node pred = null, curr = null, succ = null;
2
      boolean[] marked = {false};
      boolean snip:
5 retry: while (true) {
        pred = head; curr = pred.next.getReference();
6
        while (true) {
7
          succ = curr.next.get(marked);
8
          while (marked[0]) {
9
            snip = pred.next.compareAndSet(curr. succ. false, false);
10
            if (!snip) continue retry:
            curr = succ; succ = curr.next.get(marked);
12
13
          if (curr.key >= key)
14
            return new Window(pred, curr);
15
          pred = curr; curr = succ;
16
        } // end while(true) on line 7
17
      ł
18
```
Concurrent Set with Nonblocking Synchronization

```
public boolean add(T x) {
1
     int key = x.hashcode();
2
     while (true) {
3
       Window w = find(head, key);
4
        Node pred = w.pred, curr = w.curr;
5
        if (curr.key == key) return false;
       else {
7
          Node node = new Node(x);
8
          node.next = new AtomicMarkableReference(curr. false):
9
          if (pred.next.compareAndSet(curr, node, false, false))
10
            return true;
11
12
13
14
```

Concurrent Set with Nonblocking Synchronization

```
public boolean remove(T x) {
  int key = x.hashcode();
  boolean snip;
 while (true) {
   Window w = find(head, key);
    Node pred = w.pred, curr = w.curr;
    if (curr.kev != kev) return false:
   else {
      Node succ = curr.next.getReference();
      snip = curr.next.compareAndSet(succ, succ, false, true);
      if (!snip) continue:
      pred.next.compareAndSet(curr. succ. false, false);
      return true;
```

1

2

3

4

5

6

8

9

10

11

Concurrent Set with Nonblocking Synchronization

```
public boolean contains(T x) {
    int key = x.hashcode();
    Node curr = head;
    while (curr.key < key) {
        curr = curr.next.getReference();
    }
    return curr.key == key && !curr.next.isMarked();
}</pre>
```

ABA Problem

Thread 1 is executing deq(a)



Assume that deleted nodes can be reused

Thread 1 is executing deq(a)



Thread 1 sees head points to a, but gets delayed before executing the CAS in deq(a)

Thread 1 is executing deq(a)

Other threads execute deq(a, b, c, d), then execute enq(a)



Thread 1 is executing deq(a)

Other threads execute deq(a, b, c, d), then execute enq(a)



Thread 1's CAS succeeds, incorrectly setting head to the recycled node b

Avoiding ABA Problem

- Common workaround is to add extra "tag" to the memory address being compared
 - ▶ Tag can be a counter that tracks the number of updates to the reference
 - Can steal lower order bits of memory address or use a separate tag field if 128-bit CAS is available
- LL/SC does not suffer from the ABA problem because it checks whether a value has changed in between the interval, rather than comparing the value itself

Concurrent Queues

Bounded Partial Queue



Enqueue and dequeue operations are at the two ends — allows for concurrent modifications

Bounded Partial Queue

Given these requirements, what do we need to have a correct concurrent implementation?

Given these requirements, what do we need to have a correct concurrent implementation?

- Lock for mutual exclusion of concurrent enqueues
- Lock for mutual exclusion of concurrent dequeues
- Condition variable to indicate queue is empty
- Condition variable to indicate queue is full
- Atomic variable to track the current size

Bounded Partial Queue: enq()

```
public void eng(T x) {
     boolean wakeDeg = false;
2
     Node e = new Node(x):
3
     engLock.lock();
     trv {
5
       while (size.get() == MAX_CAPACITY)
6
         notFull.await():
       // Linearization point
8
       tail.next = e:
9
       // Update the tail pointer
10
       tail = e:
11
       if (size.getAndIncrement() == 0)
         wakeDeg = true:
     } finallv {
14
       engLock.unlock();
15
16
```

```
if (wakeDeg) {
    deaLock.lock():
    trv {
      notEmpty.signalAll();
    } finally {
      deqLock.unlock();
} // end eng()
```

17

18

10

20

21

22

23

24

25

26

27

28

30

31

Bounded Partial Queue: deq()

```
public void deg() {
                                                   if (wakeEng) {
                                              15
     boolean wakeEng = false;
                                                      engLock.lock():
                                              16
     T result:
                                                      trv {
3
                                              17
     degLock.lock();
                                                        notFull.signalAll();
                                              18
                                                      } finally {
     trv {
5
                                              10
                                                        engLock.unlock():
       while (head.next == null)
6
                                              20
         notEmpty.await():
                                              21
       result = head.next.value:
8
                                              22
       head = head.next:
                                                   return result:
                                              23
9
       if (size.getAndDecrement() ==
                                                 }
10
                                              24
           MAX CAPACITY)
                                              25
         wakeEng = true:
11
                                              26
     } finally {
12
                                              27
       deaLock.unlock():
                                              28
14
                                              29
```

Evaluating the Bounded Queue

- Need to ensure correct interleaving of concurrent calls to enq() and deq()
 - ► Special cases: Queue has zero or one element (size can become negative temporarily)
- Shared updates to the size variable could be a bottleneck
 - Can we do something about it?

Unbounded Total Queue

enq() always enqueues an item, will ignore OOM errors
deq() returns an error if the queue is empty

Simpler conditions, no need for condition variables and no need to track the size

```
public void eng(T x) {
     Node e = new Node(x):
2
                                                  2
     engLock.lock();
                                                  3
     trv {
4
                                                  4
       tail.next = e:
5
                                                  5
       tail = e:
6
                                                  6
     } finally {
                                                  7
        engLock.unlock();
8
                                                  8
9
                                                  9
10
                                                 10
                                                 11
                                                 12
                                                 13
                                                    }
13
```

```
public T deg() {
  T result:
  deqLock.lock();
  trv {
    if (head.next == null)
      return null:
    result = head.next.value:
    head = head.next:
  } finally {
    deqLock.unlock();
  return result:
```

Unbounded Lock-free Queue

```
public class Node {
    public T val;
    public AtomicReference<Node> next;
3
    public Node(T value) {
       this.val = val;
       next = AtomicReference<Node>(null);
6
8
9
  public class LockFreeQueue<T> {
10
    AtomicReference<Node> head. tail:
    public LockFreeQueue() {
       Node node = new Node(null);
       head = new AtomicReference(node):
14
       end = new AtomicReference(node);
15
     }
17
     . . .
18
```

Unbounded Lock-free Queue: enq()

```
public void eng(T x) {
  Node node = new Node(x):
  while (true) {
    Node last = tail.get(), next = last.next.get();
    if (last == tail.get()) {
      if (next == null) {
        if (last.next.compareAndSet(next. node)) {
          tail.compareAndSet(last, node);
          return;
      } else { // if (next == null)
        // Indicates a concurrent engueuer. tail not vet updated
        tail.compareAndSet(last. next):
    } // end while (true)
```

2

3

4

5

6

7

8

9 10

11

13 14

Unbounded Lock-free Queue: enq()

```
public void eng(T x) {
     Node node = new Node(x):
2
     while (true) {
3
       Node last = tail.get(). next = last.next.get();
4
       if (last == tail.get()) {
         if (next == null) {
6
           if (last.next.compareAndSet(next, node)) {
7
             tail.compareAndSet(last, node);
8
             return:
9
10
         } else { // if (next == null)
11
           // Indicates a concurrent engueuer. tail not vet updated
           tail.compareAndSet(last, next);
13
14
         // end while (true)
15
                                                 Where is the linearization
16
                                                  point?
17
```

Unbounded Lock-free Queue: dea()

```
public void deq(T x) throws EmptyException {
    while (true) {
       Node first = head.get(), last = tail.get();
       Node next = first.next.get();
       if (first == head.get()) {
         if (first == last) {
           if (next == null)
             throw new EmptyException();
          // tail is lagging behind head
           tail.compareAndSet(last, next);
         } else { // There are valid nodes
           T val = next.value;
           if (head.compareAndSet(first. next))
             return val:
       } // end if (first == head.get())
    } // end while (true)
18 }
```

2

3

4

5

6

7

8

9

10

12

13

14 15

16

Ensure that tail remains valid







Ensure that tail remains valid



Concurrent Hash Sets

Closed addressing Each table entry refers to a set of items, called a bucket. Closed addressing is also known as chaining.

Open addressing Each table entry maps to a single item. Open addressing requires a deterministic probing scheme to search for free slots.

Let *l* be the number of probing attempts, *c* be the capacity of the hash set, and s(k, l) the *l*-th element in the probe sequence where s(k, o) = h(k).

Linear probing Probing sequence is $s(k, l) = (h(k) + l) \mod c$. Cache efficient but suffers from clustering.

Quadratic probing Probing sequence is $s(k, l) = (h(k) + l^2) \mod c$. Incurs more cache misses but avoids primary clustering.

Chaotic probing Probing sequence is $s(k, l) = (h(k) + l \cdot g(k)) \mod c$ where g(k) is a second hash function. Incurs more cache misses but avoids primary clustering. Chaotic probing is also known as double hashing.

Swarnendu Biswas (IIT Kanpur)

Hash Set with Closed Addressing: Abstract Base Class

```
public abstract class BaseHashSet<T> {
  protected List<T>[] table;
  protected int setSize;
  public BaseHashSet(int capacity) {
    setSize = 0;
    table = (List<T>[]) new List[capacity];
    for (int i = 0: i < capacity: i++)</pre>
      table[i] = new ArrayList<T>();
  public boolean contains(T x) {
    acquire(x);
    trv {
      int myBucket = x.hashCode() % table.length:
      return table[mvBucket].contains(x):
    } finally {
      release(x):
```

1

2

3

4

6

7

8

10

12

14

16

Hash Set with Closed Addressing: Abstract Base Class

```
public boolean add(T x) {
19
        boolean result = false:
20
        acquire(x);
21
        trv {
22
          int myBucket = x.hashCode() % table.length;
23
          result = table[mvBucket].add(x):
24
          setSize = result ? setSize + 1 : setSize;
25
        } finally {
26
          release(x);
27
28
        if (policy()) // When to resize the hash set?
29
          resize():
30
        return result;
31
32
    } // end class BaseHashSet<T>
33
```

Policies: average bucket size exceeds a fixed threshold, more than 1/4 of the buckets exceed a bucket threshold, or if any single bucket exceeds a global threshold

Swarnendu Biswas (IIT Kanpur)

Hash Set with Closed Addressing: Coarse-Grained Locking

```
public class CoarseHashSet<T> extends BaseHashSet<T>{
  final Lock lock:
  CoarseHashSet(int capacity) {
    super(capacity);
    lock = new ReentrantLock();
  public final void acquire(T x) {
    lock.lock():
  public void release(T x) {
    lock.unlock():
  public boolean policv() {
    // Average size of a bucket is > 4
    return setSize / table.length > 4;
```

1

2

3

4

6

7

8 9

10

11

14

Hash Set with Closed Addressing: Coarse-Grained Locking

```
public void resize() {
  int oldCapacity = table.length;
  lock.lock();
  trv {
    if (oldCapacity != table.length)
      return; // Someone beat us to it
    int newCapacity = 2 * oldCapacity;
    List<T>[] oldTable = table;
    table = (List<T>[]) new List[newCapacity];
    for (int i = 0: i < newCapacity: i++)</pre>
      table[i] = new ArravList<T>():
    for (List<T> bucket : oldTable)
      for (T x : bucket)
        table[x.hashCode() % table.length].add(x);
  } finally {
    lock.unlock():
```

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

Hash Set with Closed Addressing: Fine-Grained Locking

- How would you design a concurrent hash table with fine-grained synchronization?
- Remember that a dynamically-sized hash table may require resizing

Hash Set with Closed Addressing: Fine-Grained Locking

- How would you design a concurrent hash table with fine-grained synchronization?
- Remember that a dynamically-sized hash table may require resizing

- Associating a lock with every hash table entry would incur large space overhead
- Accesses to the hash table are expected to be more distributed when the table is large implying low contention
- Resizing the lock array is complicated
 - Several locks may be in use, there can be concurrent resizes, need to avoid deadlocks
 - Delays threads that have initiated a concurrent operation
 - The hash index computed before resizing may no longer be valid after resizing

Hash Set with Closed Addressing: Striped Locking



- Each lock protects N/L buckets
- Allows more concurrency than coarse-grained lock

Hash Set with Closed Addressing: Striped Locking

```
public class StripedHashSet<T> extends BaseHashSet<T>{
  final ReentrantLock[] locks;
  public StripedHashSet(int capacity) {
    super(capacity);
    // Number of locks is initially same as the length of the table. The table
    // can dynamically grow, but not locks. Increasing the number of locks is
    // challenging.
    locks = new Lock[capacity];
    for (int j = 0; j < locks.length; j++)</pre>
      locks[j] = new ReentrantLock();
  public final void acquire(T x) {
    locks[x.hashCode() % locks.length].lock():
  public void release(T x) {
    locks[x.hashCode() % locks.length].unlock();
```

1

2

3

4

5

6

7

8

9

10 11

12

13 14

15

Hash Set with Closed Addressing: Striped Locking

```
public void resize() {
  int oldCapacity = table.length;
  for (Lock lock : locks)
    lock.lock():
  try {
    if (oldCapacity != table.length) return; // Someone beat us to it
    int newCapacity = 2 * oldCapacity:
    List<T>[] oldTable = table;
    table = (List<T>[]) new List[newCapacity];
    for (int i = 0: i < newCapacity: i++)</pre>
      table[i] = new ArrayList<T>();
    for (List<T> bucket : oldTable)
      for (T x : bucket)
        table[x.hashCode() % table.length].add(x);
  } finally {
    for (Lock lock : locks)
      lock.unlock():
```

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

Hash Set with Closed Addressing: Refinable Hash Set

Allow resizing number of locks

2

3

4

5

6

7

8

9

10

11

13

14 15

16

17

```
public class RefinableHashSet<T> extends BaseHashSet<T> {
 // Identifies the owner thread who is resizing, and the boolean flag is set to true.
 // Used for mutual exclusion with other mutation methods (e.g., add()).
 AtomicMarkableReference<Thread> owner:
 volatile ReentrantLock[] locks;
 public RefinableHashSet(int capacity) {
    super(capacity);
    locks = new ReentrantLock[capacity];
    for (int i = 0; i < capacity; i++)</pre>
     locks[i] = new ReentrantLock();
    owner = new AtomicMarkableReference<Thread>(null, false);
 public void release(T x) {
    locks[x.hashCode() % locks.length].unlock():
 protected void quiesce() { // Visit each lock and wait until it is free
    for (ReentrantLock lock : locks)
     while (lock.isLocked()) {}
```

Hash Set with Closed Addressing: Refinable Hash Set

```
public void acquire(T x) {
22
      boolean[] mark = {true};
23
      Thread me = Thread.currentThread();
24
      Thread who:
      while (true) {
26
        do { // Wait while some other thread is the owner
          who = owner.get(mark):
        } while (mark[0] && who != me);
        ReentrantLock[] oldLocks = locks;
        ReentrantLock oldLocks = oldLocks[x.hashCode() % oldLocks.length];
31
        oldLock.lock():
        // Check again to see if the locks array has been resized in the meantime
        who = owner.get(mark):
34
        // locks array has not changed, mark is not set or mark is set and I am the owner
35
        if ((!mark[0] || who == me) && locks == oldLocks) {
          return;
        } else {
          oldLock.unlock():
        }
42
```

25

27

28

29

30

32

33

36

37

38

39
Hash Set with Closed Addressing: Refinable Hash Set

```
public void resize() {
42
        boolean[] mark = {false};
43
        int oldCapacity = table.length, newCapacity = 2 * oldCapacity;
44
        Thread me = Thread.currentThread();
45
        if (owner.compareAndSet(null, me, false, true)) { // Try to make yourself the owner
46
          trv {
47
            if (table.length != oldCapacity) return; // Someone else resized first
48
            quiesce():
49
            List<T>[] oldTable = table:
50
            table = (List<T>[]) new List[newCapacity];
51
            for (int i = 0: i < newCapacity: i++)</pre>
52
               table[i] = new ArravList<T>():
53
            locks = new ReentrantLock[newCapacity];
54
            for (int j = 0; j < locks.length; j++)</pre>
55
               locks[j] = new ReentrantLock();
56
            initializeFrom(oldTable); // Copy old data
57
          } finally {
58
            owner.set(null. false):
59
60
61
62
```

Hash Set with Closed Addressing: Lock-free Hash Set

Challenging to design a correct algorithm with synchronization primitives like CAS

- Difficult to "stop-the-world" while resizing the bucket
- Not enough to make individual buckets lock-free
- Resizing the table requires atomically moving entries from old buckets to new buckets
- If the table doubles in capacity, then items in the old bucket must be distributed between two new buckets
- However, a CAS operates on only one memory location

Hash Set with Open Addressing: Cuckoo Hashing

Cuckoo hashing is an open addressing scheme where collisions are resolved by displacing any earlier item occupying the same slot with a newly added item

- Assume a hash set of size N = 2k, and two tables each of size k (denoted by table[0] and table[1]
- Two independent hash functions h_0 and h_1 map the keys to $0, \ldots, k-1$ providing two possible locations for each key

contains(x) Tests whether either table[0][$h_0(x)$] or table[1][$h_1(x)$] is equal to x

remove(x) Checks whether x is in either table[0][$h_0(x)$] or table[1][$h_1(x)$] and removes it if found

- add(x) Repeatedly displace conflicting items until every key has a slot
 - May not find an empty slot if the table is full or the sequence of displacements form a cycle
 - May need to resize the hash table, choose new hash functions, and restart the add operation after a THRESHOLD of successive displacements is reached

Sequential Cuckoo Hashing: add()





Concurrent Cuckoo Hashing

Challenge is in the possibly long sequence of swap operations during add()

Break up each method call into a sequence of phases, where each phase adds, removes, or displaces a single item x

- Hash table is organized as a 2D table of probe sets
 - > Probe set is a constant-sized set of items with the same hash code
 - Each probe set holds at most PROBE_SIZE items
- Implementation tries to ensure that when the set is quiescent (i.e., no method calls are in progress), each probe set holds no more than THRESHOLD (< PROBE_SIZE) items

Concurrent Cuckoo Hashing: Abstract Hash Table Structure

```
public abstract class PhasedCuckooHashSet<T> {
1
      volatile int capacity:
2
      volatile List<T>[][] table:
3
      public PhasedCuckooHashSet(int size) {
4
        capacity = size;
        table = (List<T>[][]) new java.util.ArrayList[2][capacity];
6
        for (int i = 0; i < 2; i++) {</pre>
7
          for (int j = 0; j < capacity; j++) {</pre>
8
            table[i][j] = new ArrayList<T>(PROBE SIZE);
9
10
11
12
```

Concurrent Cuckoo Hashing: Addition and Relocation



Example of relocation

Example of adding to the Hash Set

Swarnendu Biswas (IIT Kanpur)

Concurrent Cuckoo Hashing: Remove

```
public boolean remove(T x) {
14
      acquire(x);
15
      trv {
16
        List<T> seto = table[o][hasho(x) % capacity]:
        if (seto.contains(x)) {
18
           set0.remove(x);
19
           return true;
20
        } else {
21
           List<T> set1 = table[1][hash1(x) % capacity];
22
           if (set1.contains(x)) {
23
             set1.remove(x):
24
             return true;
25
26
27
        return false:
28
      } finally {
29
        release(x):
30
31
32
```

Concurrent Cuckoo Hashing: Addition

```
public boolean add(T x) {
                                                                    } else {
33
      T y = null;
                                                                     mustResize = true:
34
                                                           52
      acquire(x);
35
      int ho = hasho(x) % capacity;
                                                                 } finallv {
36
                                                           54
      int h_1 = hash_1(x) \% capacity:
                                                                   release(x):
37
                                                          55
      int i = -1, h = -1;
38
                                                          56
      boolean mustResize = false:
                                                                 // Either resize or relocate
39
      try {
                                                                 if (mustResize) {
40
                                                          58
        if (contains(x)) return false;
                                                                   resize():
41
                                                          50
        List<T> seto = table[o][ho]:
                                                                   add(x):
42
                                                          60
        List<T> set1 = table[1][h1];
                                                                 } else if (!relocate(i, h)) {
43
                                                           61
        if (seto.size() < THRESHOLD) {</pre>
                                                                   // Item already added
44
                                                          62
           seto.add(x): return true:
                                                                   resize();
45
                                                           63
         } else if (set1.size() < THRESHOLD) {</pre>
46
                                                          64
           set1.add(x): return true:
47
                                                          65
                                                                 return true:
         } else if (seto.size() < PROBE SIZE) {</pre>
48
                                                          66
           set0.add(x); i = 0; h = h0; // Relocate
                                                          67
49
         } else if (set1.size() < PROBE SIZE) {</pre>
                                                          68
50
           set1.add(x); i = 1; h = h1; // Relocate
                                                          69
51
```

Concurrent Cuckoo Hashing: Relocation

```
protected boolean relocate(int i, int hi) {
70
      int hj = 0, j = 1 - i;
71
      for (int round=0: round < LIMIT: round++) {</pre>
72
        List<T> iSet = table[i][hi]; // Probe set to shrink
73
        T v = iSet.get(o); // Oldest element is the candidate
74
        switch (i) {
75
          case 0: hj = hash1(y) % capacity; break;
76
          case 1: hi = hasho(v) % capacity: break:
77
78
        acquire(y);
79
        List<T> jSet = table[j][hj]; // Other probe set
80
        trv {
81
          if (iSet.remove(y)) {
82
             if (jSet.size() < THRESHOLD) {</pre>
83
               iSet.add(y); return true;
84
             } else if (jSet.size() < PROBE SIZE) {</pre>
85
               jSet.add(y); // Relocate other probe set
86
               i = 1 - i; hi = hj; j = 1 - j;
87
             } else { // jSet is full, trigger resize
88
               iSet.add(y); return false;
89
90
```

Concurrent Cuckoo Hashing: Relocation



Concurrent Cuckoo Hashing using Striped Locking

```
public class StripedCuckooHashSet<T> extends PhasedCuckooHashSet<T>{
  final ReentrantLock[][] lock;
  public StripedCuckooHashSet(int capacity) {
    super(capacity);
    lock = new ReentrantLock[2][capacity];
    for (int i = 0; i < 2; i++) {</pre>
      for (int j = 0; j < capacity; j++)</pre>
        lock[i][j] = new ReentrantLock():
  public final void acquire(T x) {
    lock[0][hash0(x) % lock[0].length].lock();
    lock[1][hash1(x) % lock[1].length].lock();
  public final void release(T x) {
    lock[0][hasho(x) % lock[0].length].unlock();
    lock[1][hash1(x) % lock[1].length].unlock();
```

2

3

4

5

6

7

8 9 10

12

13 14

15

16

17 18

Concurrent Cuckoo Hashing using Striped Locking

```
public void resize() {
  int oldCapacity = capacity;
  for (Lock aLock : lock[0]) { aLock.lock(): }
  trv {
    if (capacity != oldCapacity) { return; }
    List<T>[][] oldTable = table;
    capacity = 2 * capacity;
    table = (List<T>[][]) new List[2][capacity];
    for (List<T>[] row : table) {
      for (int i = 0: i < row.length: i++) {</pre>
        row[i] = new ArrayList<T>(PROBE SIZE):
      }
    for (List<T>[] row : oldTable) {
      for (List<T> set : row) {
        for (T z : set) \{ add(z); \}
  } finally {
    for (Lock aLock : lock[0]) { aLock.unlock(); }
  }
```

19

20

21

22

24

25

26

27

28

29

30

31

32

33 34

35

36

37 38 39

References

- M. Herlihy et al. The Art of Multiprocessor Programming. Chapters 3, 9, 10, 11, 13, 2nd edition, Morgan Kaufmann.
- Remzi H. Arpaci-Dusseau and Andrea C. Arpaci-Dusseau. Operating Systems: Three Easy Pieces. Chapter 29, Online.



- Mark Moir and Nir Shavit. Concurrent Data Structures In Handbook of Data Structures and Applications, 2nd edition, Chapman and Hall/CRC Press, 2004.
- A. Castañeda et al. A Linearizability-based Hierarchy for Concurrent Specifications. In Communications of the ACM, volume 66, issue 1, 2022, pp. 86–97.



A. Castañeda and S. Rajsbaum. Recent Advances on Principles of Concurrent Data Structures. In Communications of the ACM, volume 67, issue 8, 2024, pp. 45–46.