

# CS 610: Concurrent Data Structures

CPU's and GPU's

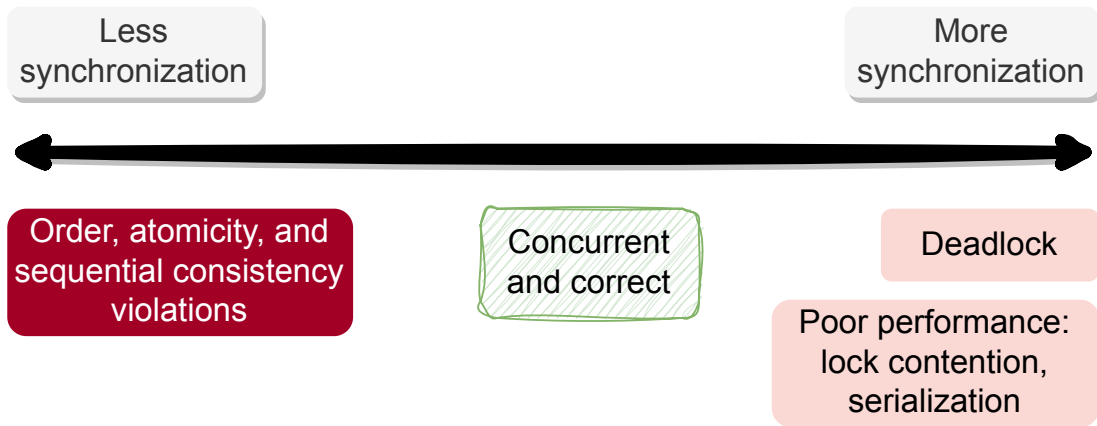
**Swarnendu Biswas**

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

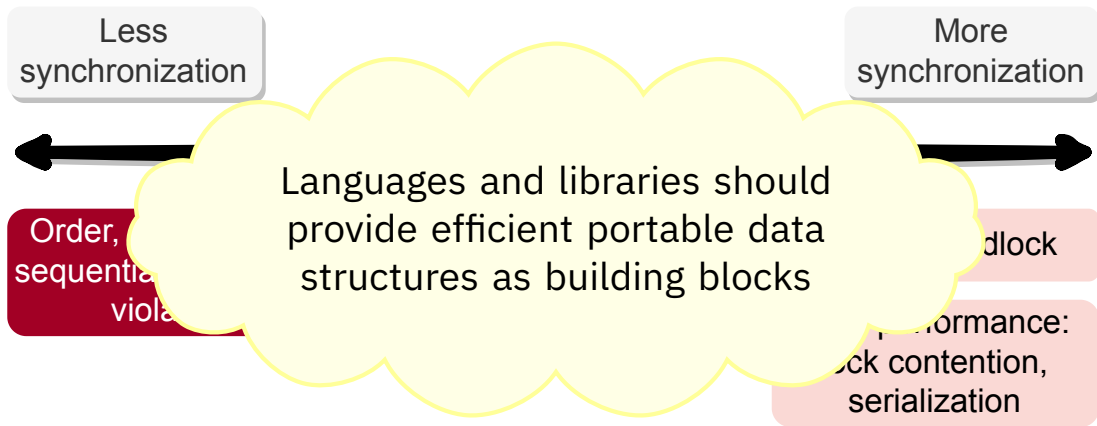
Sem 2024-25-I



# Challenges with Concurrent Programming



# Challenges with Concurrent Programming



# Designing a Concurrent Set Data Structure

# Designing a Concurrent Set

```
1 public interface Set<T> {  
2     boolean add(T x);  
3     boolean remove(T x);  
4     boolean contains(T x);  
5 }
```

## 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)

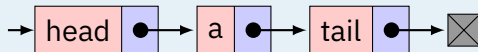
returns true if and only if x is present in the set

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

# Designing a Concurrent Set Using Linked Lists

```
1 class Node {  
2     T data;  
3     int key;  
4     Node next;  
5 }
```

Two sentinel nodes head and tail



## 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

```
1 public class UnsafeList<T> {  
2     private Node head;  
3     public UnsafeList() {  
4         head = new Node(Integer.MIN_VALUE);  
5         head.next = new Node(Integer.  
6             MAX_VALUE);  
7     }  
8 }
```

```
1 public boolean add(T x) {  
2     int key = x.hashCode();  
3     Node pred = head;  
4     Node curr = pred.next;  
5     while (curr.key < key) {  
6         pred = curr; curr = curr.next;  
7     }  
8     if (key == curr.key) {  
9         return false;  
10    } else {  
11        Node node = new Node(x);  
12        node.next = curr;  
13        prev.next = node;  
14        return true;  
15    }  
16 }
```

# A Thread-Unsafe Set Data Structure

```
1  public boolean remove(T x) {
2      int key = x.hashCode();
3      Node pred = head;
4      Node curr = pred.next;
5      while (curr.key < key) {
6          pred = curr;
7          curr = curr.next;
8      }
9      if (key == curr.key) {
10         pred.next = curr.next;
11         return true;
12     } else {
13         return false;
14     }
15 }
```

```
1  public boolean contains(T x) {
2      int key = x.hashCode();
3      Node pred = head;
4      Node curr = pred.next;
5      while (curr.key < key) {
6          pred = curr;
7          curr = curr.next;
8      }
9      if (key == curr.key) {
10         return true;
11     } else {
12         return false;
13     }
14 }
15 }
```



# A Thread-Unsafe Set Data Structure

```
1 public boolean remove(T x) {  
2     int key = x.hashCode();  
3     Node pred = head;  
4     Node curr = pred;  
5     while (curr != null) {  
6         if (curr.key == key) {  
7             curr = curr.next;  
8             return true;  
9         }  
10        curr = curr.next;  
11    }  
12    return false;  
13 }  
14  
15 }
```

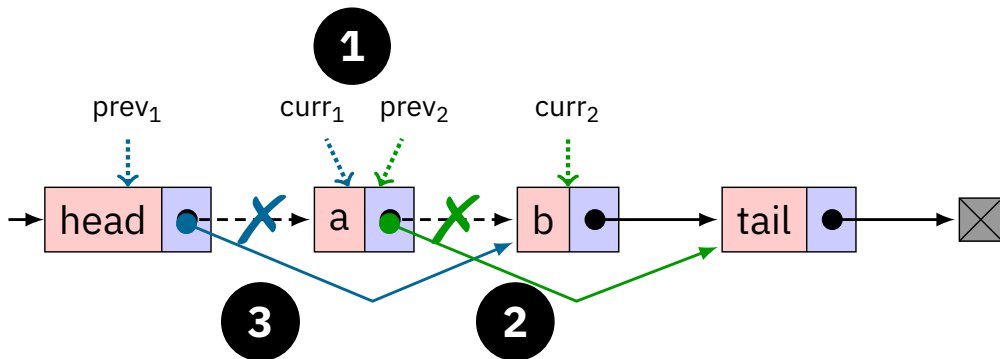
```
1 public boolean contains(T x) {  
2     int key = x.hashCode();  
3     Node pred = head;  
4     while (pred != null) {  
5         if (pred.key == key) {  
6             return true;  
7         }  
8         pred = pred.next;  
9     }  
10    return false;  
11 }  
12  
13  
14  
15 }
```

Can you give an example  
to show that remove()  
is not thread-safe?

# Thread-Unsafe Set: Incorrect remove()

Thread 1 is executing remove(a)

Thread 2 is executing remove(b)



# Concurrent Set with Coarse-Grained Synchronization

```
1 public class CoarseList<T> {  
2     private Node head;  
3     private Lock lock = new ReentrantLock();  
4  
5     public CoarseList() {  
6         head = new Node(Integer.MIN_VALUE);  
7         head.next = new Node(Integer.MAX_VALUE);  
8     }  
9     ...  
10  
11    public boolean add(T x) {  
12        Node pred, curr  
13        int key = x.hashCode();  
14        lock.lock();
```

```
15        try {  
16            pred = head;  
17            curr = pred.next;  
18            while (curr.key < key) {  
19                pred = curr;  
20                curr = curr.next;  
21            }  
22            if (key == curr.key) {  
23                return false;  
24            } else {  
25                Node node = new Node(x);  
26                node.next = curr;  
27                prev.next = node;  
28                return true;  
29            }  
30        } finally {  
31            lock.unlock();  
32        }  
33    }
```

# Concurrent Set with Coarse-Grained Synchronization

```
34 public boolean remove(T x) {
35     Node pred, curr;
36     int key = x.hashCode();
37     lock.lock();
38     try {
39         pred = head;
40         curr = pred.next;
41         while (curr.key < key) {
42             pred = curr;
43             curr = curr.next;
44         }
45         if (key == curr.key) {
46             pred.next = curr.next;
47             return true;
48         } else {
49             return false;
50         }
51     } finally {
52         lock.unlock();
53     }
54 }
```

```
55 public boolean contains(T x) {
56     Node curr;
57     int key = x.hashCode();
58     boolean found = false;
59     lock.lock();
60     try {
61         curr = head.next;
62         while (curr.key < key) {
63             curr = curr.next;
64         }
65         if (key == curr.key) {
66             found = true;
67         }
68     } finally {
69         lock.unlock();
70     }
71     return found;
72 }
73 }
```

# Concurrent Set with Fine-Grained Synchronization

Add a lock object to each list node

```
1  class Node {  
2      T data;  
3      int key;  
4      Node key;  
5      Lock lock;  
6  }
```

## Possible interleaving

### Thread 1

```
1  curr.lock.lock();  
2  next = curr.next;  
3  curr.lock.unlock();  
4  
5  next.lock.lock();
```

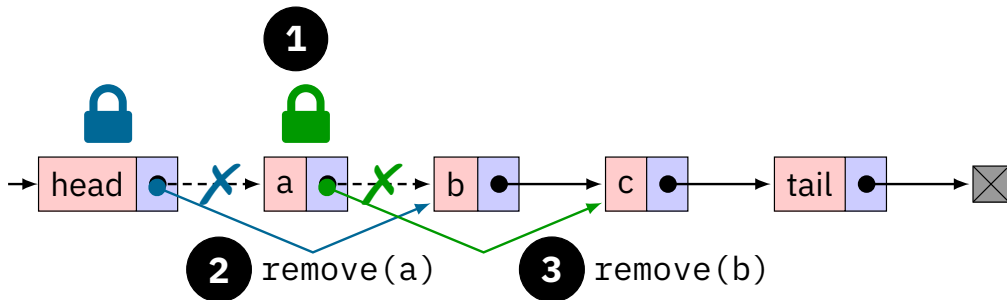
### Thread 2

```
1  
2  
3  
4  // Remove next from list  
5
```

# Is Locking One Node Sufficient?

Thread 1 is executing `remove(a)`

Thread 2 is executing `remove(b)`



# Concurrent Set with Fine-Grained Synchronization

```
1  public boolean add(T x) {  
2      int key = x.hashCode();  
3      head.lock();  
4      Node pred = head;  
5      try {  
6          Node curr = pred.next;  
7          curr.lock();  
8          try {  
9              while (curr.key < key) {  
10                 pred.unlock();  
11                 pred = curr;  
12                 curr = curr.next;  
13                 curr.lock();  
14             }  
15         }
```

```
16         if (key == curr.key) {  
17             return false;  
18         } else {  
19             Node node = new Node(x);  
20             node.next = curr;  
21             pred.next = node;  
22             return true;  
23         }  
24     } finally {  
25         curr.unlock();  
26     }  
27 } finally {  
28     pred.unlock();  
29 }  
30 }
```

# Concurrent Set with Fine-Grained Synchronization

```
1  public boolean remove(T x) {  
2      int key = x.hashCode();  
3      head.lock();  
4      Node pred = null, curr = null;  
5      try {  
6          pred = head; curr = pred.next;  
7          curr.lock();  
8          try {  
9              while (curr.key < key) {  
10                 pred.unlock();  
11                 pred = curr;  
12                 curr = curr.next;  
13                 curr.lock();  
14             }  
15         }  
16     }  
17 }
```

```
15         if (key == curr.key) {  
16             pred.next = curr.next;  
17             return true;  
18         } else {  
19             return false;  
20         }  
21     } finally {  
22         curr.unlock();  
23     }  
24 } finally {  
25     pred.unlock();  
26 }  
27 }  
28 }
```



# 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

# 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  $q.\text{enq}(z)$  changes the state of the queue to  $Q \bullet z$ , 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 structure

Situation:

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?

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

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

# Reasoning about Correctness of Concurrent Data Structures

- 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

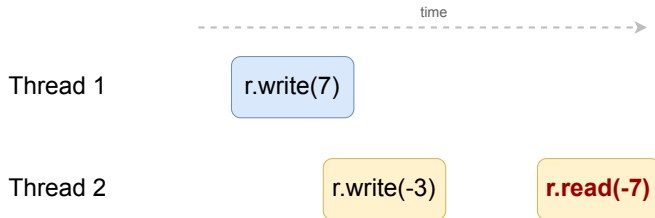
# Definitions

- 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  $P$ .

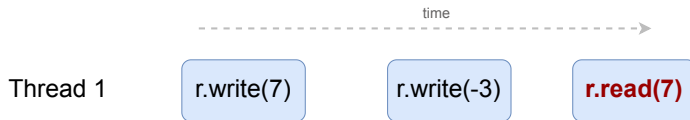


# Correct Behavior Expected From a Concurrent Execution

(i) Method calls should appear to happen one-at-a-time in sequential order



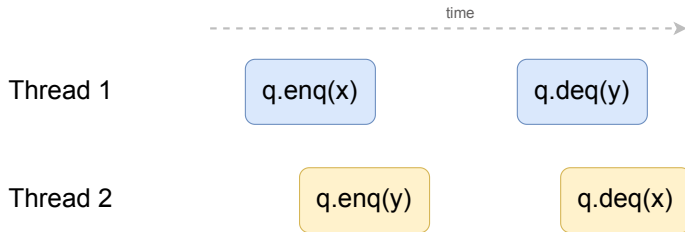
(ii) Method calls should appear to take effect in program order



# 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.

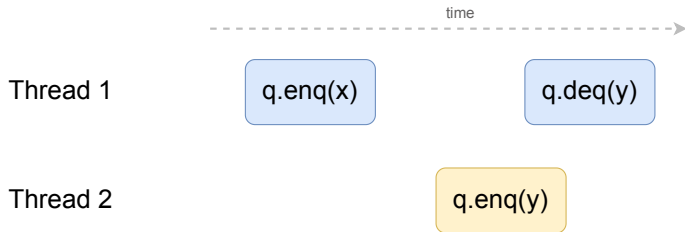
Consider the following operations on a FIFO queue  $q$ , where  $x$  and  $y$  are objects.



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

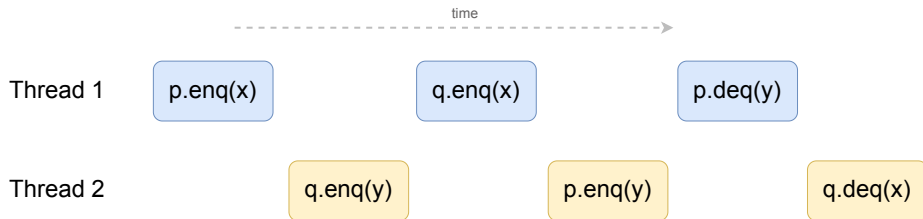
# SC can Violate Real-Time Order

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

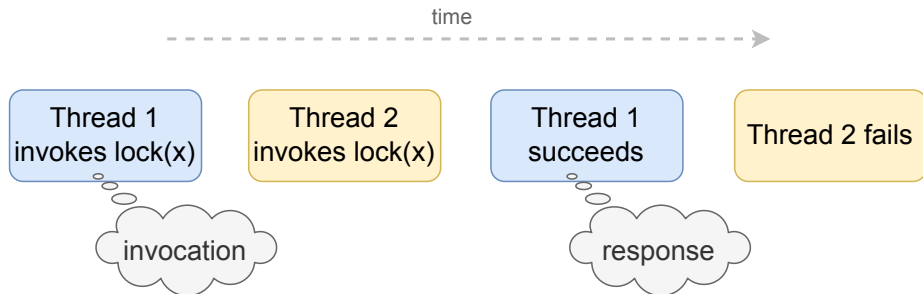
## Linearizability has two requirements

- (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

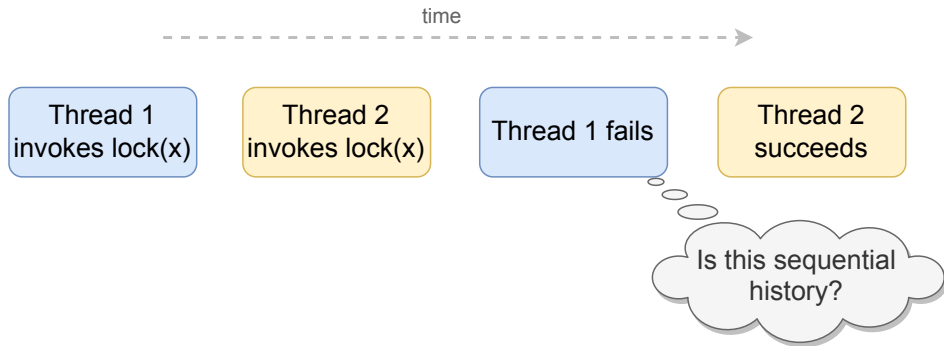
# Understanding Linearizability

- 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



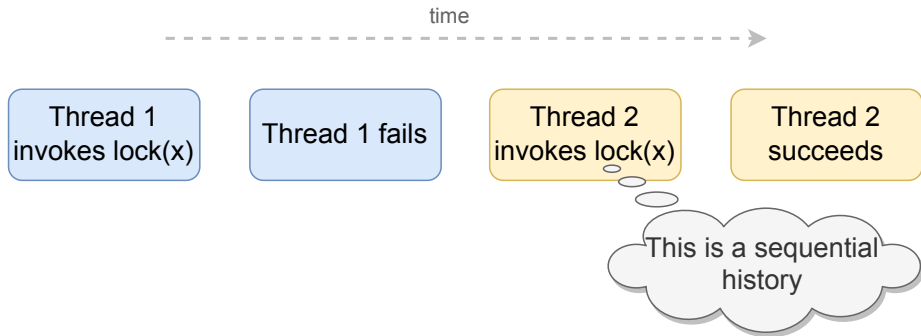
# Sequential History

- 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





# Understanding Linearizability

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

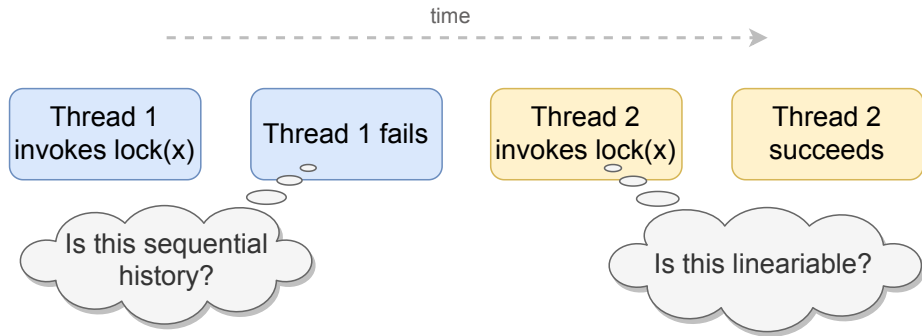
Consider a concurrent history (set of method calls)  $H$  and a valid sequential history  $S$ . The history  $H$  is linearizable if:

- 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  $S$

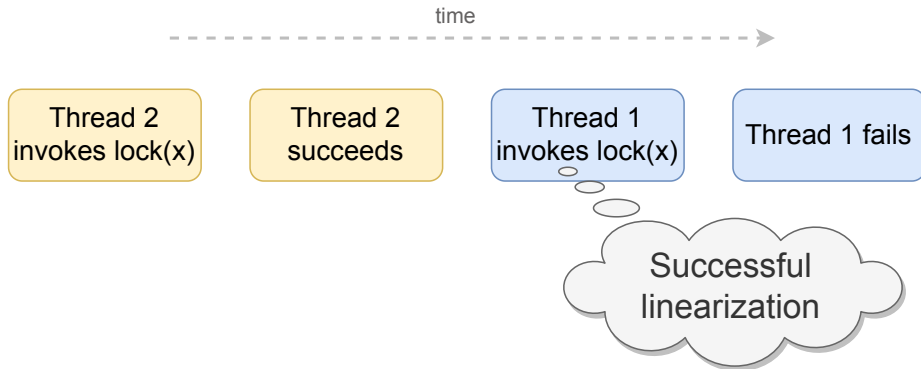
# Linearizability in Simpler Words

- 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

# 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

# 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

# Concurrent Set with Optimistic Synchronization

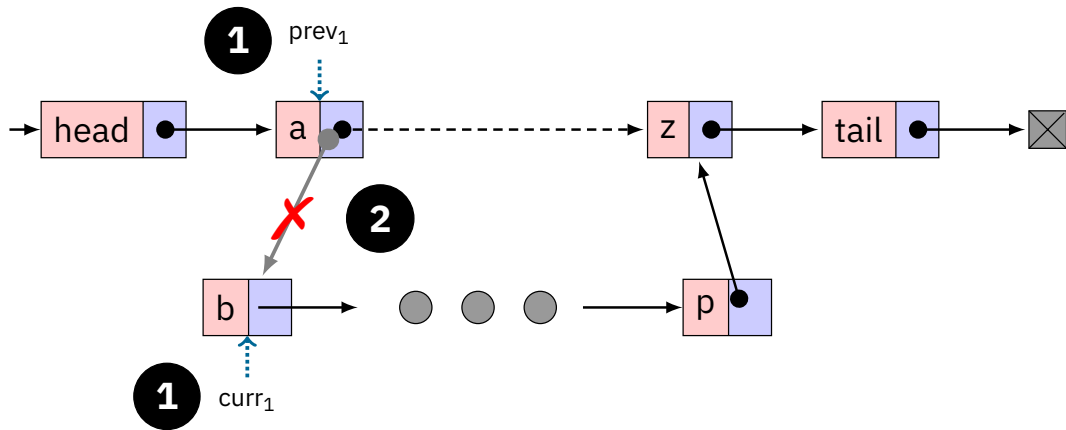
```
1  public boolean add(T x) {
2      int key = x.hashCode();
3      while (true) {
4          Node pred = head;
5          Node curr = pred.next;
6          while (curr.key < key) {
7              pred = curr;
8              curr = curr.next;
9          }
10         pred.lock(); curr.lock();
```

```
17     try {
18         if (validate(pred, curr)) {
19             if (curr.key == key) {
20                 return false;
21             } else {
22                 Node node = new Node(x);
23                 node.next = curr;
24                 prev.next = node;
25                 return true;
26             }
27         }
28     } finally {
29         curr.unlock(); pred.unlock();
30     }
31 }
32 }
```

# 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`

```
1  boolean validate(Node prev, Node curr) {  
2      Node node = head;  
3      while (node.key <= prev.key) {  
4          if (node == prev)  
5              return prev.next == curr;  
6          node = node.next;  
7      }  
8      return false;  
9  }
```

# Concurrent Set with Optimistic Synchronization

```
1  public boolean remove(T x) {  
2      int key = x.hashCode();  
3      while (true) {  
4          Node pred = head;  
5          Node curr = pred.next;  
6          while (curr.key < key) {  
7              pred = curr;  
8              curr = curr.next;  
9          }  
10         pred.lock(); curr.lock();  
11  
12  
13  
14
```

```
15         try {  
16             if (validate(pred, curr)) {  
17                 if (curr.key == key) {  
18                     pred.next = curr.next;  
19                     return true;  
20                 } else {  
21                     return false;  
22                 }  
23             }  
24             finally {  
25                 curr.unlock(); pred.unlock();  
26             }  
27         }  
28     }
```

# Concurrent Set with Optimistic Synchronization

```
1  public boolean contains(T x) {  
2      int key = x.hashCode();  
3      while (true) {  
4          Node pred = head;  
5          Node curr = pred.next;  
6          while (curr.key < key) {  
7              pred = curr;  
8              curr = curr.next;  
9          }  
10         pred.lock();  
11         curr.lock();
```

```
12         try {  
13             if (validate(pred, curr)) {  
14                 return curr.key == key;  
15             }  
16         } finally {  
17             curr.unlock();  
18             pred.unlock();  
19         }  
20     }  
21 }  
22
```

# Concurrent Set with Optimistic Synchronization

## Are there problems with the optimistic-synchronization-based Set design?

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



# 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

```
1 public boolean add(T x) {
2     int key = x.hashCode();
3     while (true) {
4         Node pred = head;
5         Node curr = pred.next;
6         while (curr.key < key) {
7             pred = curr;
8             curr = curr.next;
9         }
10        pred.lock();
11        try {
12            curr.lock();
13            try {
14                if (validate(pred, curr)) {
15                    if (curr.key == key) {
16                        return false;
```

```
17            } else {
18                Node node = new Node(x);
19                node.next = curr;
20                pred.next = node;
21                return true;
22            }
23        }
24        } finally {
25            curr.unlock(); }
26    }
27    } finally {
28        pred.unlock();
29    }
30    }
31 }
32 }
```

# How to Validate?

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

```
1 boolean validate(Node prev, Node curr) {  
2     return !prev.marked && !curr.marked && prev.next == curr;  
3 }
```

# Concurrent Set with Lazy Synchronization

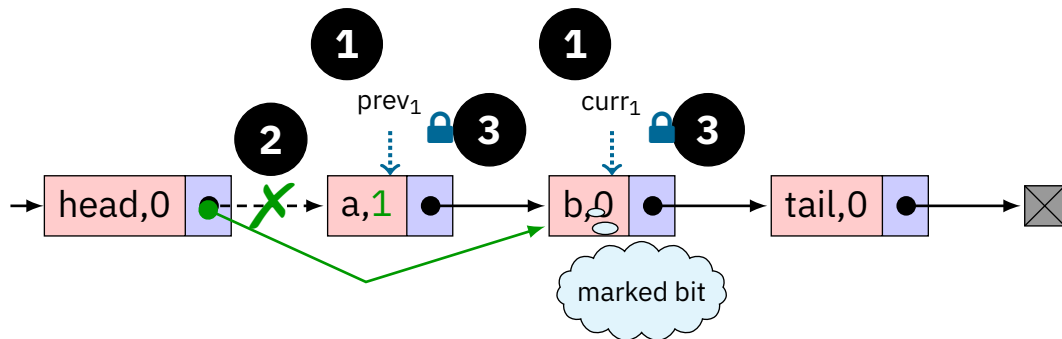
```
1  public boolean remove(T x) {
2      int key = x.hashCode();
3      while (true) {
4          Node pred = head;
5          Node curr = pred.next;
6          while (curr.key < key) {
7              pred = curr;
8              curr = curr.next;
9          }
10         pred.lock();
11         try {
12             curr.lock();
13             try {
14                 if (validate(pred, curr)) {
15                     if (curr.key != key) {
16                         return false;
```

```
17         } else {
18             // Logical deletion
19             curr.marked = true;
20             // Physical deletion
21             pred.next = curr.next;
22             return true;
23         }
24     }
25     } finally {
26         curr.unlock(); }
27     }
28     } finally {
29         pred.unlock();
30     }
31 }
32 }
```

# Detecting Conflicts: Scenario 1

Thread 1 is executing remove (b)

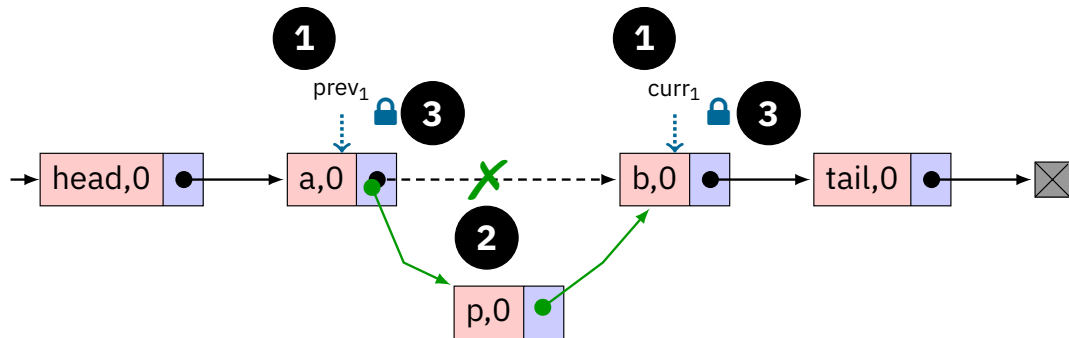
Thread 2 is executing remove (a)



# Detecting Conflicts: Scenario 2

Thread 1 is executing `remove(b)`

Thread 2 is executing `add(p)`



# Concurrent Set with Lazy Synchronization

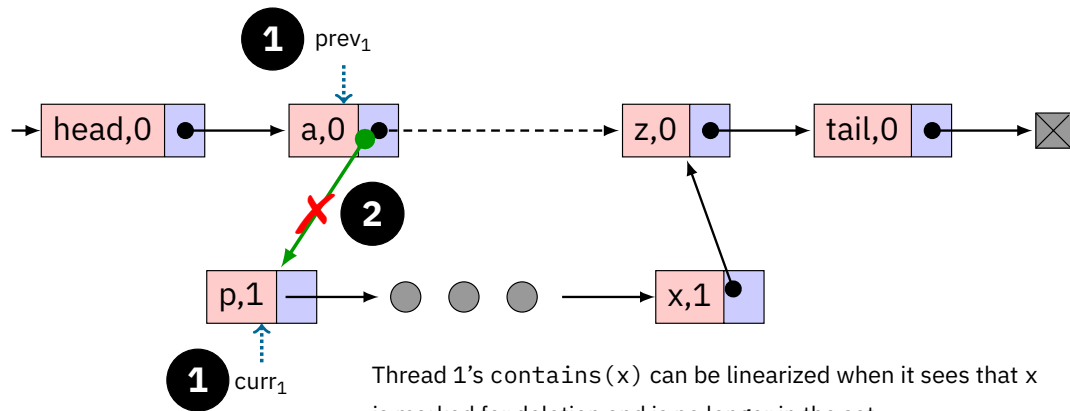
```
1 public boolean contains(T x) {  
2     int key = x.hashCode();  
3     Node curr = head;  
4     while (curr.key < key) {  
5         curr = curr.next;  
6     }  
7     return curr.key == key && !curr.marked;  
8 }
```



# Unsuccessful contains(): Scenario 1

Thread 1 is executing contains(x)

Thread 2 is executing remove(p...x)

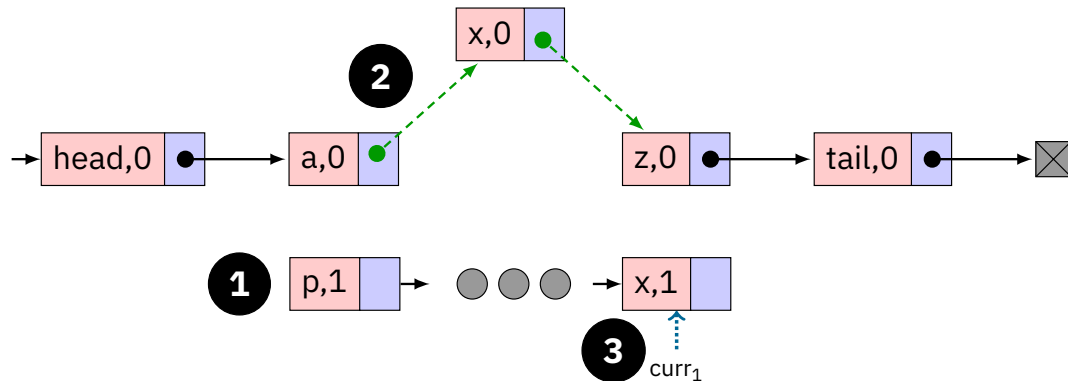




## Unsuccessful contains(): Scenario 2

Thread 1 is executing contains(x)

Thread 2 is executing add(x)



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

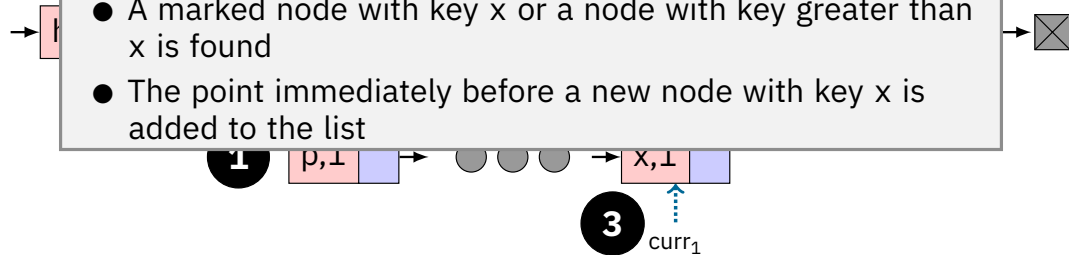
## Unsuccessful contains(): Scenario 2

Thread 1 is executing contains(x)

Thread 2 is executing add(x)

Linearize an unsuccessful contains(x) at the earlier of the following two points:

- A marked node with key x or a node with key greater than x is found
- The point immediately before a new node with key x is added to the list



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

# Nonblocking Synchronization

## Why do we need nonblocking designs?

- 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
- Use of locks can lead to deadlocks, livelocks, and priority inversion

**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

```
1 bool CAS(word* loc, word oldval, word newval) {  
2     atomic { // Code block will execute atomically  
3         res := (*loc == oldval);  
4         if (res)  
5             *loc := newval;  
6         return res;  
7     }  
8 }
```

# 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

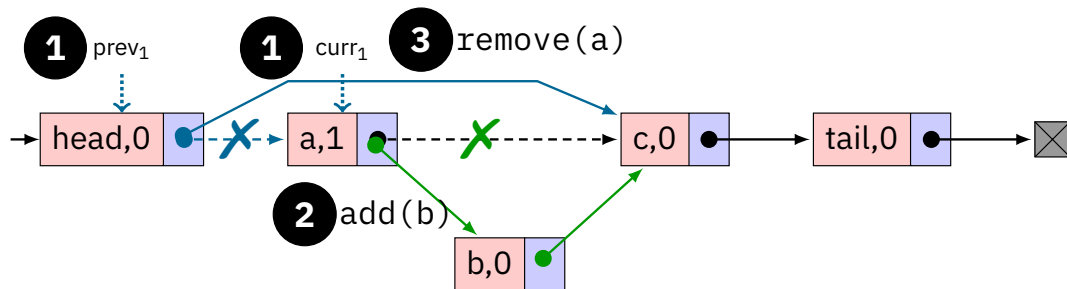
```
1      xor %ecx, %ecx /*ecx=0*/
2      inc %ecx /*ecx=1*/
3  RETRY: xor %eax, %eax /*eax=0*/
4          lock cmpxchg %ecx, &lock
5          jnz RETRY
6          ret
7
```

```
1  void spinLock(lock* lk) {
2      // flag attribute is set when the
3      // lock is acquired
4      while (CAS(&lk->flag, 0, 1) == 1) {
5          // Keep spinning
6      }
7  }
```

# Nonblocking Synchronization with CAS

Thread 1 is executing remove (a)

Thread 2 is executing add (b)

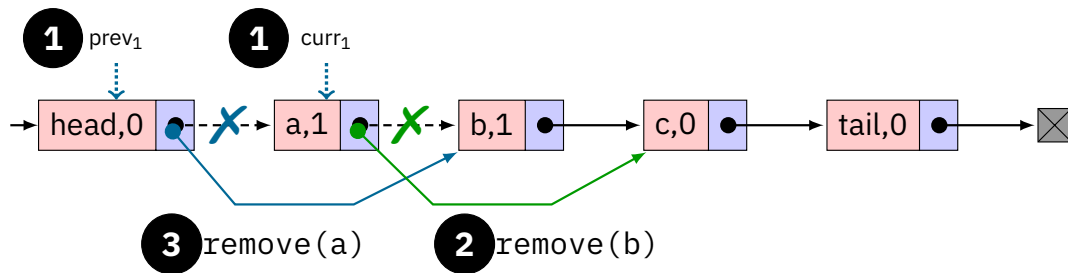


a is deleted but b is not added to the list

# Nonblocking Synchronization with CAS

Thread 1 is executing remove (a)

Thread 2 is executing remove (b)



a is deleted but b is not deleted from the list



# 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

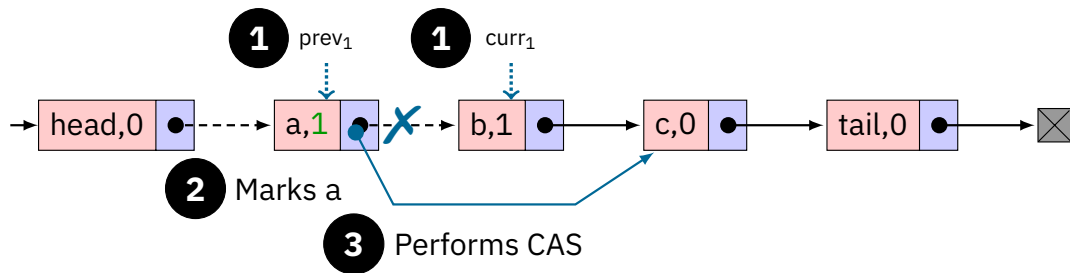


Why?

# Challenge in Traversing Marked Nodes

Thread 1 is executing `remove(b)`

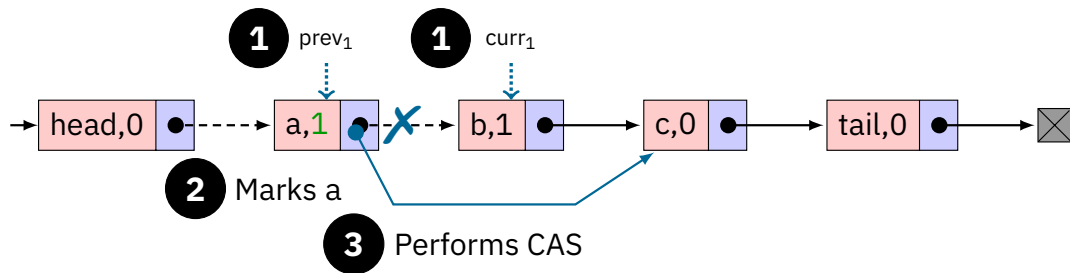
Thread 2 marks `a`



# 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`

# Helper Method

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`

```
1 class Window {  
2     public Node pred, curr;  
3     Window(Node myPred, Node myCurr) {  
4         pred = myPred; curr = myCurr;  
5     }  
6 }
```

# Helper Method

```
1 public Window find(Node head, int key) {
2     Node pred = null, curr = null, succ = null;
3     boolean[] marked = {false};
4     boolean snip;
5 retry: while (true) {
6     pred = head; curr = pred.next.getReference();
7     while (true) {
8         succ = curr.next.get(marked);
9         while (marked[0]) {
10             snip = pred.next.compareAndSet(curr, succ, false, false);
11             if (!snip) continue retry;
12             curr = succ; succ = curr.next.get(marked);
13         }
14         if (curr.key >= key)
15             return new Window(pred, curr);
16         pred = curr; curr = succ;
17     }
18 }
```

# Concurrent Set with Nonblocking Synchronization

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

# Concurrent Set with Nonblocking Synchronization

```
1 public boolean remove(T x) {  
2     int key = x.hashCode();  
3     boolean snip;  
4     while (true) {  
5         Window w = find(head, key);  
6         Node pred = w.pred, curr = w.curr;  
7         if (curr.key != key) return false;  
8         else {  
9             Node succ = curr.next.getReference();  
10            snip = curr.next.compareAndSet(succ, succ, false, true);  
11            if (!snip) continue;  
12            pred.next.compareAndSet(curr, succ, false, false);  
13            return true;  
14        }  
15    }  
16 }
```

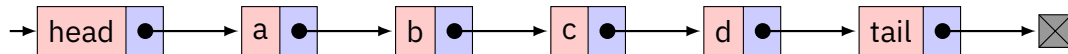


# Concurrent Set with Nonblocking Synchronization

```
1 public boolean contains(T x) {  
2     int key = x.hashCode();  
3     Node curr = head;  
4     while (curr.key < key) {  
5         curr = curr.next.getReference();  
6     }  
7     return curr.key == key && !curr.next.isMarked();  
8 }
```

# Lock-Free Programming and ABA Problem

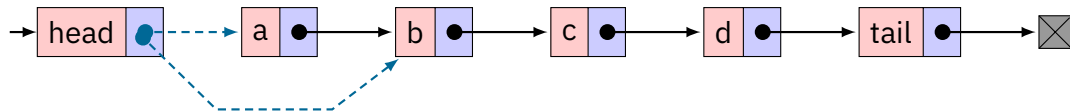
Thread 1 is executing `deq(a)`



Assume deleted nodes can be reused

# Lock-Free Programming and ABA Problem

Thread 1 is executing `deq(a)`

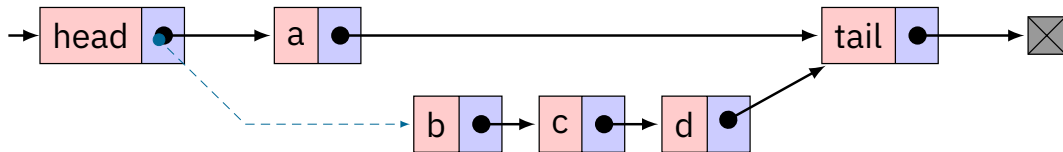


Thread 1 sees head points to a, but gets delayed while executing `deq(a)`

# Lock-Free Programming and ABA Problem

Thread 1 is executing `deq(a)`

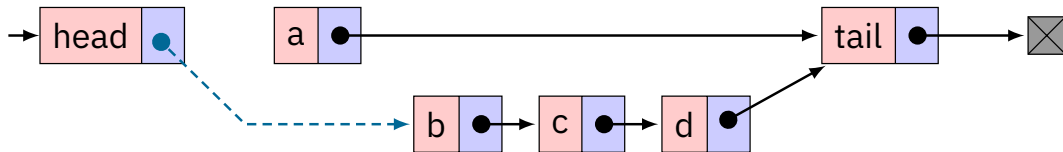
Other threads execute `deq(a, b, c, d)`, then execute `enq(a)`



# Lock-Free Programming and ABA Problem

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

# Concurrent Hash Sets for CPU

# 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  $s(k, 0) = h(k)$ .

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

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

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



# Hash Set with Closed Addressing: Abstract Base Class

```
1 public abstract class BaseHashSet<T> {  
2     protected List<T>[] table;  
3     protected int setSize;  
4     public BaseHashSet(int capacity) {  
5         setSize = 0;  
6         table = (List<T>[]) new List[capacity];  
7         for (int i = 0; i < capacity; i++)  
8             table[i] = new ArrayList<T>();  
9     }  
10    public boolean contains(T x) {  
11        acquire(x);  
12        try {  
13            int myBucket = x.hashCode() % table.length;  
14            return table[myBucket].contains(x);  
15        } finally {  
16            release(x);  
17        }  
18    }
```

# Hash Set with Closed Addressing: Abstract Base Class

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

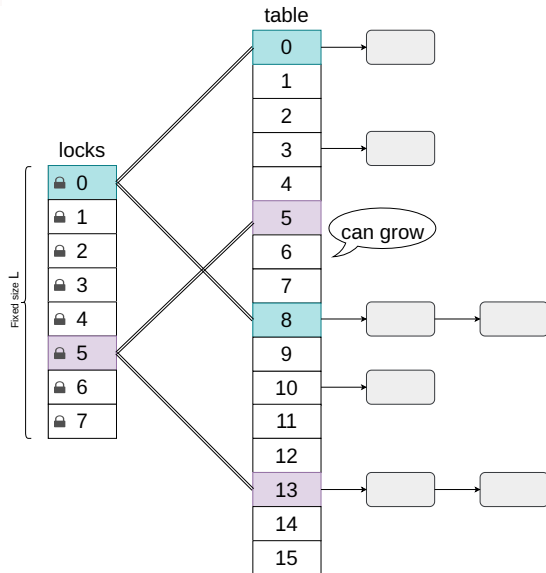
# Hash Set with Closed Addressing: Coarse-Grained Locking

```
1 public class CoarseHashSet<T> extends BaseHashSet<T>{
2     final Lock lock;
3     CoarseHashSet(int capacity) {
4         super(capacity);
5         lock = new ReentrantLock();
6     }
7     public final void acquire(T x) {
8         lock.lock();
9     }
10    public void release(T x) {
11        lock.unlock();
12    }
13    public boolean policy() {
14        // Average size of a bucket is > 4
15        return setSize / table.length > 4;
16    }
```

# Hash Set with Closed Addressing: Coarse-Grained Locking

```
17 public void resize() {  
18     int oldCapacity = table.length;  
19     lock.lock();  
20     try {  
21         if (oldCapacity != table.length)  
22             return; // someone beat us to it  
23         int newCapacity = 2 * oldCapacity;  
24         List<T>[] oldTable = table;  
25         table = (List<T>[]) new List[newCapacity];  
26         for (int i = 0; i < newCapacity; i++)  
27             table[i] = new ArrayList<T>();  
28         for (List<T> bucket : oldTable)  
29             for (T x : bucket)  
30                 table[x.hashCode() % table.length].add(x);  
31     } finally {  
32         lock.unlock();  
33     }  
34 }
```

## 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

```
1 public class StripedHashSet<T> extends BaseHashSet<T>{
2     final ReentrantLock[] locks;
3     public StripedHashSet(int capacity) {
4         super(capacity);
5         // Number of locks is initially same as the length of the table. The table
6         // can dynamically grow, but not locks. Increasing the number of locks is
7         // challenging.
8         locks = new Lock[capacity];
9         for (int j = 0; j < locks.length; j++)
10             locks[j] = new ReentrantLock();
11     }
12     public final void acquire(T x) {
13         locks[x.hashCode() % locks.length].lock();
14     }
15     public void release(T x) {
16         locks[x.hashCode() % locks.length].unlock();
17     }
```

# Hash Set with Closed Addressing: Striped Locking

```
15 public void resize() {  
16     int oldCapacity = table.length;  
17     for (Lock lock : locks)  
18         lock.lock();  
19     try {  
20         if (oldCapacity != table.length) return; // someone beat us to it  
21         int newCapacity = 2 * oldCapacity;  
22         List<T>[] oldTable = table;  
23         table = (List<T>[]) new List[newCapacity];  
24         for (int i = 0; i < newCapacity; i++)  
25             table[i] = new ArrayList<T>();  
26         for (List<T> bucket : oldTable)  
27             for (T x : bucket)  
28                 table[x.hashCode() % table.length].add(x);  
29     } finally {  
30         for (Lock lock : locks)  
31             lock.unlock();  
32     }  
33 }  
34 }
```

# Hash Set with Closed Addressing: Refinable Hash Set

```
1 // Allow resizing number of locks
2 public class RefinableHashSet<T> extends BaseHashSet<T> {
3     // Identifies the owner thread who is resizing, and the boolean is set to true.
4     // Used for mutual exclusion with other mutation methods (e.g., add()).
5     AtomicMarkableReference<Thread> owner;
6     volatile ReentrantLock[] locks;
7     public RefinableHashSet(int capacity) {
8         super(capacity);
9         locks = new ReentrantLock[capacity];
10        for (int i = 0; i < capacity; i++)
11            locks[i] = new ReentrantLock();
12        owner = new AtomicMarkableReference<Thread>(null, false);
13    }
14    public void release(T x) {
15        locks[x.hashCode() % locks.length].unlock();
16    }
17    protected void quiesce() { // Visit each lock and wait until it is free
18        for (ReentrantLock lock : locks)
19            while (lock.isLocked()) {}
20    }
```



# Hash Set with Closed Addressing: Refinable Hash Set

```
22 public void acquire(T x) {
23     boolean[] mark = {true};
24     Thread me = Thread.currentThread();
25     Thread who;
26     while (true) {
27         do { // Wait while some other thread is the owner
28             who = owner.get(mark);
29         } while (mark[0] && who != me);
30         ReentrantLock[] oldLocks = locks;
31         ReentrantLock oldLock = oldLocks[x.hashCode() % oldLocks.length];
32         oldLock.lock();
33         // Check again to see if the locks array has been resized in the meantime
34         who = owner.get(mark);
35         // locks array has not changed, mark is not set or mark is set and I am the owner
36         if ((!mark[0] || who == me) && locks == oldLocks) {
37             return;
38         } else {
39             oldLock.unlock();
40         }
41     }
42 }
```

# Hash Set with Closed Addressing: Refinable Hash Set

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

# Hash Set with Closed Addressing: Lock-free Hash Set

## Challenging to design a correct algorithm with synchronization primitives like CAS

- 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
- CAS operates only on one memory location

# Hash Set with Open Addressing: Cuckoo Hashing

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

- Uses two hash functions which provides two possible locations for each key
- 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  $0, \dots, k - 1$

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

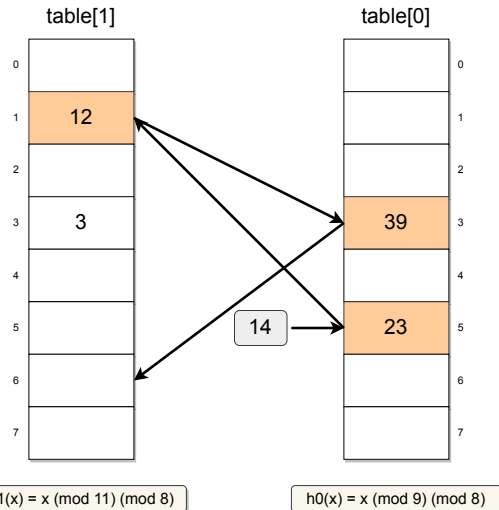
`remove(x)` Checks whether  $x$  is in either `table[0][ $h_0(x)$ ]` or `table[1][ $h_1(x)$ ]` and removes it if found

`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
- 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()

```
1 public boolean add(T x) {  
2     if (contains(x)) {  
3         return false;  
4     }  
5     for (int i = 0; i < THRESHOLD; i++) {  
6         if ((x = swap(h0(x), x)) == null) {  
7             return true;  
8         } else if ((x = swap(h1(x), x)) == null) {  
9             return true;  
10        }  
11    }  
12    resize();  
13    add(x);  
14 }
```



# Concurrent Cuckoo Hashing

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

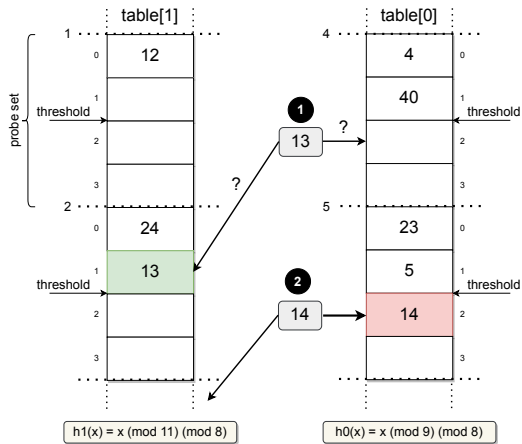
Break up each method call into a sequence of phases, where each phase adds, removes, or displaces a single item  $x$

- 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

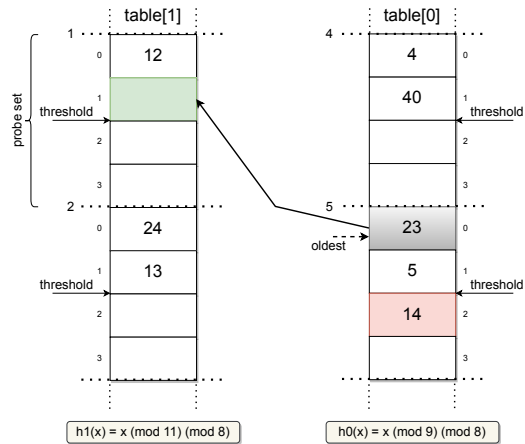
# Concurrent Cuckoo Hashing using Phases

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

# Concurrent Cuckoo Hashing using Phases



Example of adding to the Hash Set



Example of relocation



# Concurrent Cuckoo Hashing using Phases

```
13 public boolean remove(T x) {  
14     acquire(x);  
15     try {  
16         List<T> set0 = table[0][hash0(x) % capacity];  
17         if (set0.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 }
```

# Concurrent Cuckoo Hashing using Phases

```
32 public boolean add(T x) {
33     T y = null;
34     acquire(x);
35     int h0 = hash0(x) % capacity;
36     int h1 = hash1(x) % capacity;
37     int i = -1, h = -1;
38     boolean mustResize = false;
39     try {
40         if (present(x)) return false;
41         List<T> set0 = table[0][h0];
42         List<T> set1 = table[1][h1];
43         if (set0.size() < THRESHOLD) {
44             set0.add(x); return true;
45         } else if (set1.size() < THRESHOLD) {
46             set1.add(x); return true;
47         } else if (set0.size() < PROBE_SIZE) {
48             set0.add(x); i = 0; h = h0;
49         } else if (set1.size() < PROBE_SIZE) {
50             set1.add(x); i = 1; h = h1;
```

```
51         } else {
52             mustResize = true;
53         }
54     } finally {
55         release(x);
56     }
57     if (mustResize) {
58         resize();
59         add(x);
60     } else if (!relocate(i, h)) {
61         resize();
62     }
63     // x must have been present
64     return true;
65 }
```

# Concurrent Cuckoo Hashing using Phases

```
64  protected boolean relocate(int i, int
    hi) {
65      int hj = 0, j = 1 - i;
66      for (int round=0; round<LIMIT; round
          ++ ) {
67          List<T> iSet = table[i][hi];
68          T y = iSet.get(0);
69          switch (i) {
70              case 0: hj = hash1(y)%capacity;
                      break;
71              case 1: hj = hash0(y)%capacity;
                      break;
72          }
73          acquire(y);
74          List<T> jSet = table[j][hj];
75          try {
76              if (iSet.remove(y)) {
77                  if (jSet.size() < THRESHOLD) {
78                      jSet.add(y); return true;
79  
```

```
82      } else if (jSet.size() <
          PROBE_SIZE) {
83          jSet.add(y);
84          i = 1 - i; hi = hj; j = 1 -
          j;
85      } else {
86          iSet.add(y); return false;
87      }
88      } else if (iSet.size() >=
          THRESHOLD) {
89          continue;
90      } else {
91          return true;
92      }
93      } finally {
94          release(y);
95      }
96      }
97      return false;
98  }
99  }
```

# Concurrent Cuckoo Hashing using Striped Locking

```
1 public class StripedCuckooHashSet<T> extends PhasedCuckooHashSet<T>{
2     final ReentrantLock[][] lock;
3     public StripedCuckooHashSet(int capacity) {
4         super(capacity);
5         lock = new ReentrantLock[2][capacity];
6         for (int i = 0; i < 2; i++) {
7             for (int j = 0; j < capacity; j++)
8                 lock[i][j] = new ReentrantLock();
9         }
10    }
11    public final void acquire(T x) {
12        lock[0][hash0(x) % lock[0].length].lock();
13        lock[1][hash1(x) % lock[1].length].lock();
14    }
15    public final void release(T x) {
16        lock[0][hash0(x) % lock[0].length].unlock();
17        lock[1][hash1(x) % lock[1].length].unlock();
18    }
```

# Concurrent Cuckoo Hashing using Striped Locking

```
19 public void resize() {  
20     int oldCapacity = capacity;  
21     for (Lock aLock : lock[0]) {  
22         aLock.lock();  
23     }  
24     try {  
25         if (capacity != oldCapacity) {  
26             return;  
27         }  
28         List<T>[][] oldTable = table;  
29         capacity = 2 * capacity;  
30         table = (List<T>[][][]) new List[2][  
31             capacity];  
32         for (List<T>[] row : table) {  
33             for (int i = 0; i < row.length; i  
34                 ++)  
35                 row[i] = new ArrayList<T>(  
36                     PROBE_SIZE);  
37     }  
38 }
```

```
36     for (List<T>[] row : oldTable) {  
37         for (List<T> set : row) {  
38             for (T z : set) {  
39                 add(z);  
40             }  
41         }  
42     }  
43     finally {  
44         for (Lock aLock : lock[0]) {  
45             aLock.unlock();  
46         }  
47     }  
48 }  
49 }
```

# Concurrent Data Structures for GPU

# Concurrent Hashing on GPUs

## Requirements

- Need to support high throughput for concurrent accesses to the hash tables
- Need to devise nonblocking algorithms that are tuned to the GPU programming model for good performance

## Challenges in designing an efficient hash table

- Lock-based synchronization will not scale to thousands of GPU threads
- Accessing linked-list-based data structures imply making random (uncoalesced) memory accesses from threads in a warp
  - ▶ High memory bandwidth is achieved when threads in a warp access consecutive memory locations with a fixed stride
- Dynamically allocating linked list nodes for numerous GPU threads is a bottleneck

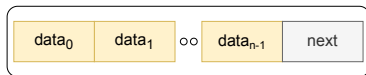
# Concurrent Hashing on GPUs with SlabHash

## SlabHash builds a concurrent hash table using slab lists instead of linked lists

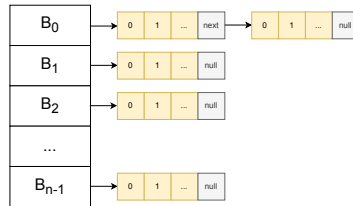
- A slab consists of multiple words of data (i.e., unordered set) and a single next pointer, thereby reducing memory overhead
- Slabs on NVIDIA GPUs can be of 32 words
- Can search a slab list node using a single ballot instruction



Linked list node



Slab list node



Slab list node



# Operations on a Slab List

- search** Search from the head of the list, load the next slab if the item cannot be found in the current slab
- insert** Start from the head of the list, use an atomic CAS to insert the new key-value pair into the first empty data slot. An unsuccessful CAS implies someone else occupied the empty slot, so restart looking for an empty slot. Load the next slab if no empty slots are found. Allocate a new slab at the end of the list if all slabs are full.
- replace** Similar to insert except that the entire slab list needs to be searched
- delete** Similar to insert except that the entire slab list needs to be searched (depending on whether duplicates are allowed)

# More about SlabHash

Provides a dynamic hash table that uses chaining for collision resolution

Universal hash function of the form  $h(k; a, b) = ((ak + b) \bmod p) \bmod B$  is used, where  $a, b$  are random arbitrary integers,  $p$  is a random prime number, and  $B$  is the number of buckets.

Threads are assigned independent tasks (i.e., keys), but work is done in parallel per-warp, called warp-cooperative work sharing (WCWS)

- One-to-one mapping maps each thread to a single key, threads in a warp process their 32 keys individually
- Advantage of WCWS is that it significantly reduces branch divergence when compared to per-thread processing
- Disadvantage of WCWS is that all threads within a warp should be active

# Concurrent Hashing on GPUs with WarpCore

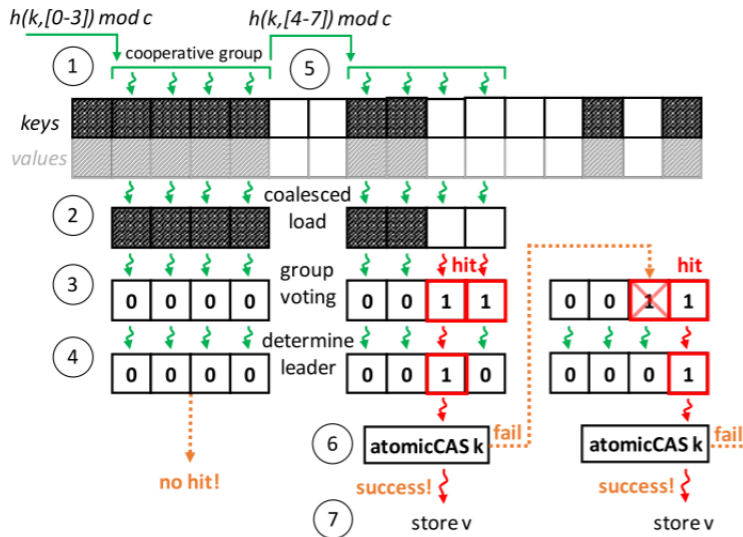
WarpCore argues that open addressing with linear probing is more amenable to the SIMT execution model on GPUs

Linked lists are cache-inefficient and lock-free insertion and deletion of nodes in linked lists is complicated because of the ABA problem



## Parallel probing scheme

- Mapping each thread to a key will lead to a different probing sequence for threads in a warp, leading to non-coalesced global memory accesses
- Can use an entire warp of 32 threads per input key  $k$ , such that each thread with lane ID  $t$  probes a different hash table position  $h(k, t) \bmod c$ , but linear probing suffers from clustering
- WarpCore uses double hashing with an inner intra-warp linear probing, double hashing determines the starting offset for linear probing

# Example of Insertion Operation with WarpCore



# References

-  M. Herlihy et al. The Art of Multiprocessor Programming. Chapters 3, 9, 10, 13, 2<sup>nd</sup> edition, Morgan and Claypool.
-  Mark Moir and Nir Shavit. Concurrent Data Structures In Handbook of Data Structures and Applications, Chapman and Hall/CRC Press, 2004.