# Concurrent Data Structures
## Concurrent Algorithms 2016

**Tudor David**

**(based on slides by Vasileios Trigonakis)**

# Data Structures (DSs)

- Constructs for <span style="color:red">efficiently storing and retrieving data</span>
  - **Different types:** lists, hash tables, trees, queues, ...
- Accessed through the <span style="color:red">DS interface</span>
  - Depends on the DS type, but always includes
  - Store an element
  - Retrieve an element
- Element
  - Set: just one value

# Concurrent Data Structures (CDSs)

- Concurrently accessed by multiple threads
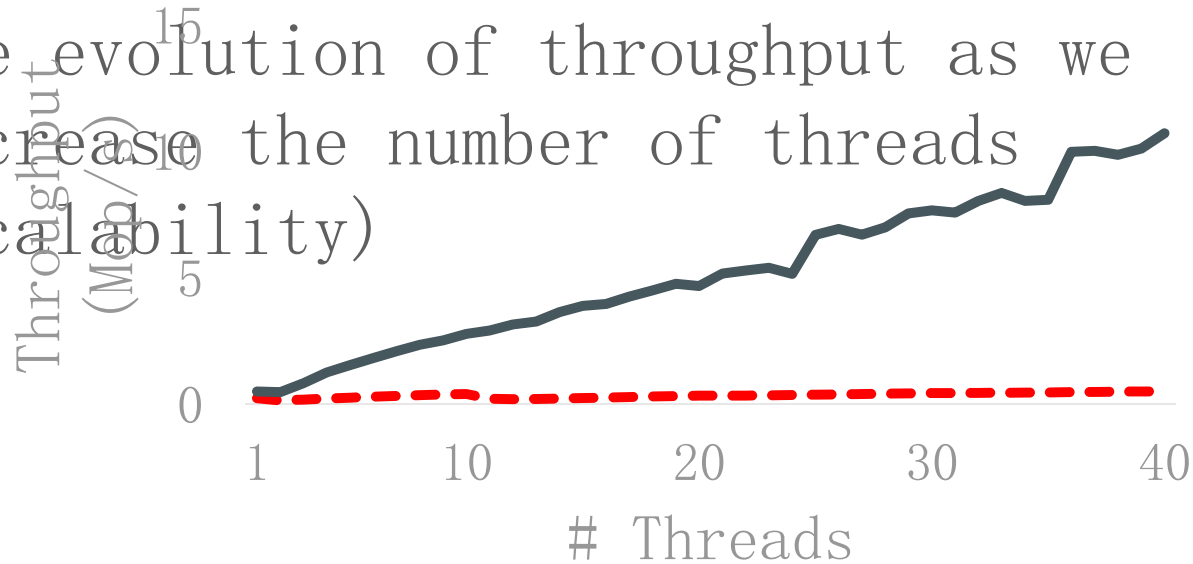  - Through the CDS interface → <span style="color:red">linearizable</span> operations!

- Really important on <span style="color:red">multi-cores</span>
- <span style="color:red">t software system</span>
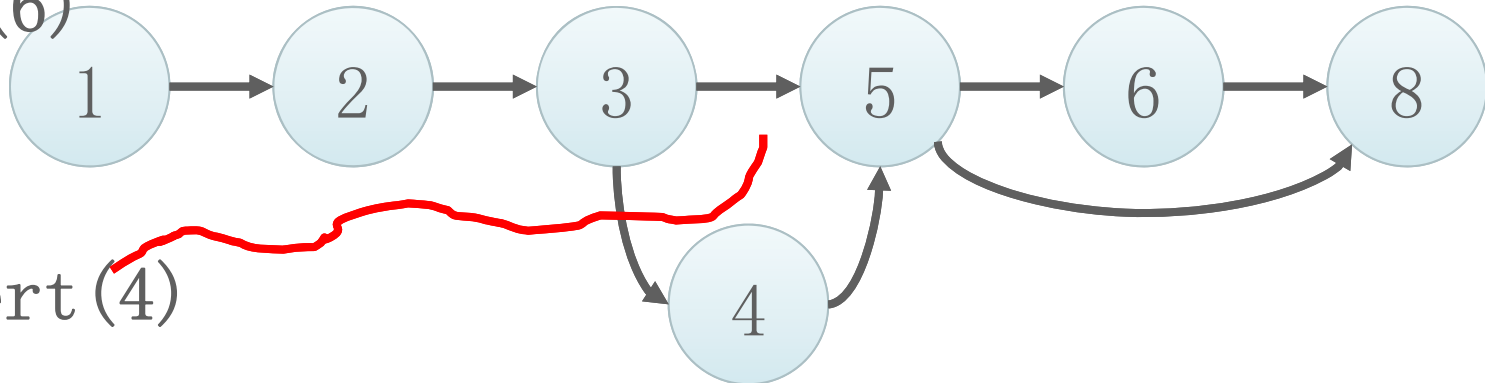
# What do we care about in practice?

- Progress of individual operations – sometimes

- More often:
  - Number of operations per second (throughput)
  - The evolution of throughput as we increase the number of threads (scalability)



Throughput
(Mop/s)

15

10

5

0

1        10        20        30        40

# Threads

# DS Example: Linked List

delete(6)

insert(4)



- A sequence of elements (nodes)

- Interface
  - search (aka contains)
  - insert
  - remove (aka delete)

```
struct node
{
  value_t value;
  struct node* next;
};
```

# Search Data Structures

- **Interface**
  1. search
  2. insert
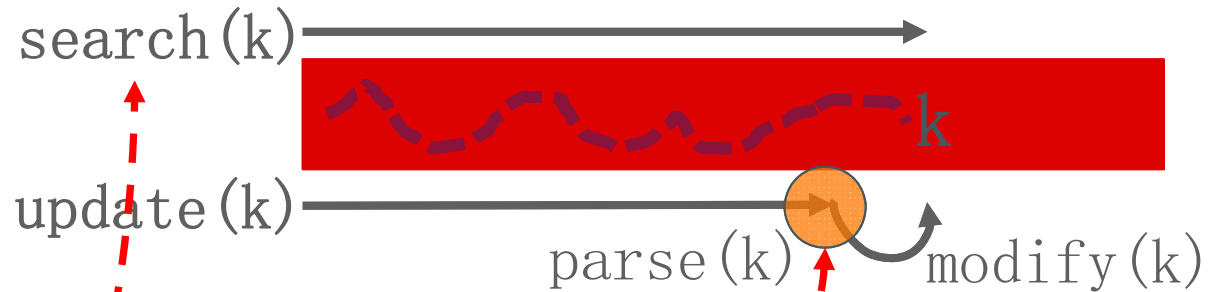     **updates**
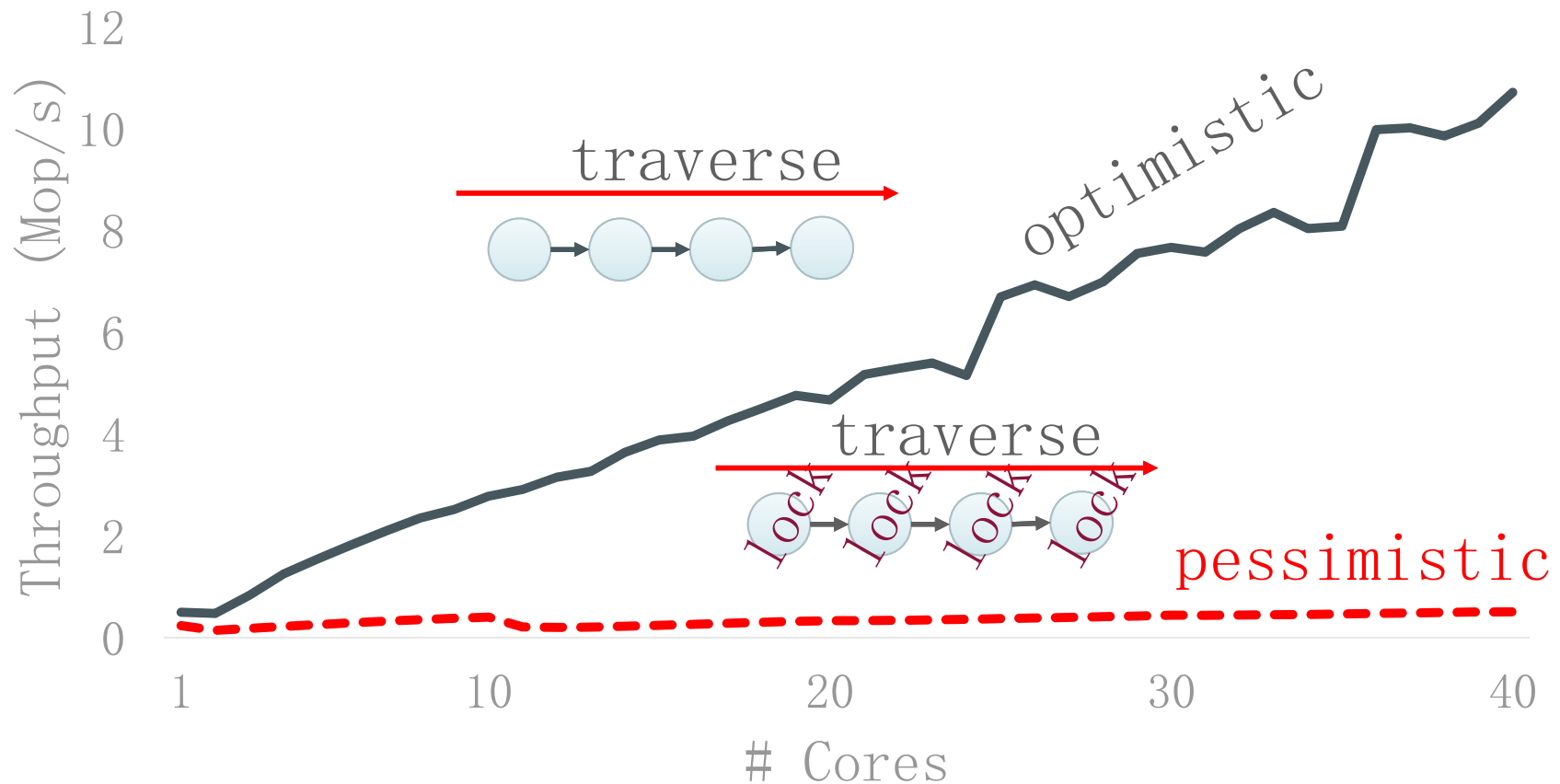  3. remove

- **Semantics**
  1. read-only
  2. read-only
  3. read-only
  4. read-write

search(k)

k

update(k)

parse(k)    modify(k)

# Optimistic vs. Pessimistic Concurrency



**(Lesson$_1$)** Optimistic concurrency is the only way to get scalability

# The Two Problems in Optimistic Concurrency

- **Concurrency Control**
  How threads synchronize their writes to the shared memory (e.g., nodes)
  - Locks
  - CAS
  - Transactional memory

- **Memory Reclamation**
  How and when threads free and reuse the shared memory (e.g., nodes)
  - Garbage collectors
  - Hazard pointers
  - RCU
  - Quiescent states

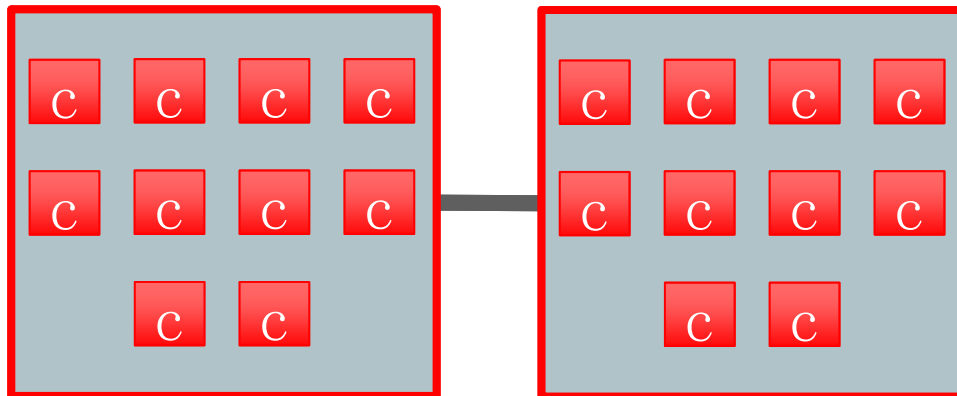# Tools for Optimistic Concurrency Control (OCC)

- RCU: slow in the presence of updates
  - (also a memory reclamation scheme)
- STM: slow in general
- HTM: not ubiquitous, not very fast (yet)

- Wait-free algorithms: slow in general
- (Optimistic) Lock-free algorithms: ☺

We either need a lock-free or an optimistic lock-based algorithm
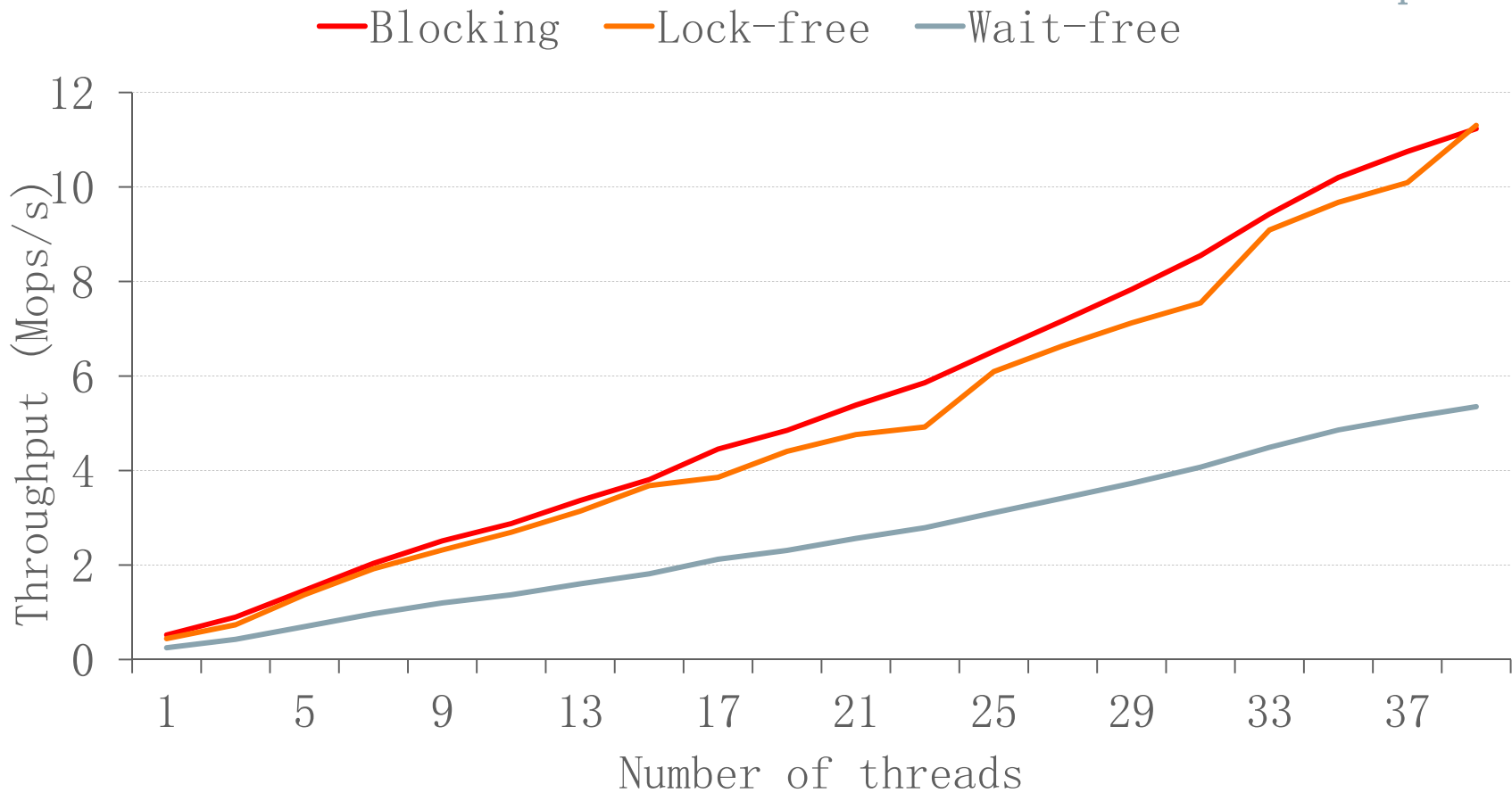
# Parenthesis: Target platform

2-socket Intel Xeon E5-2680 v2 Ivy Bridge

- 20 cores @ 2.8 GHz, 40 hyper-threads
- 25 MB LLC (per socket)
- 256GB RAM

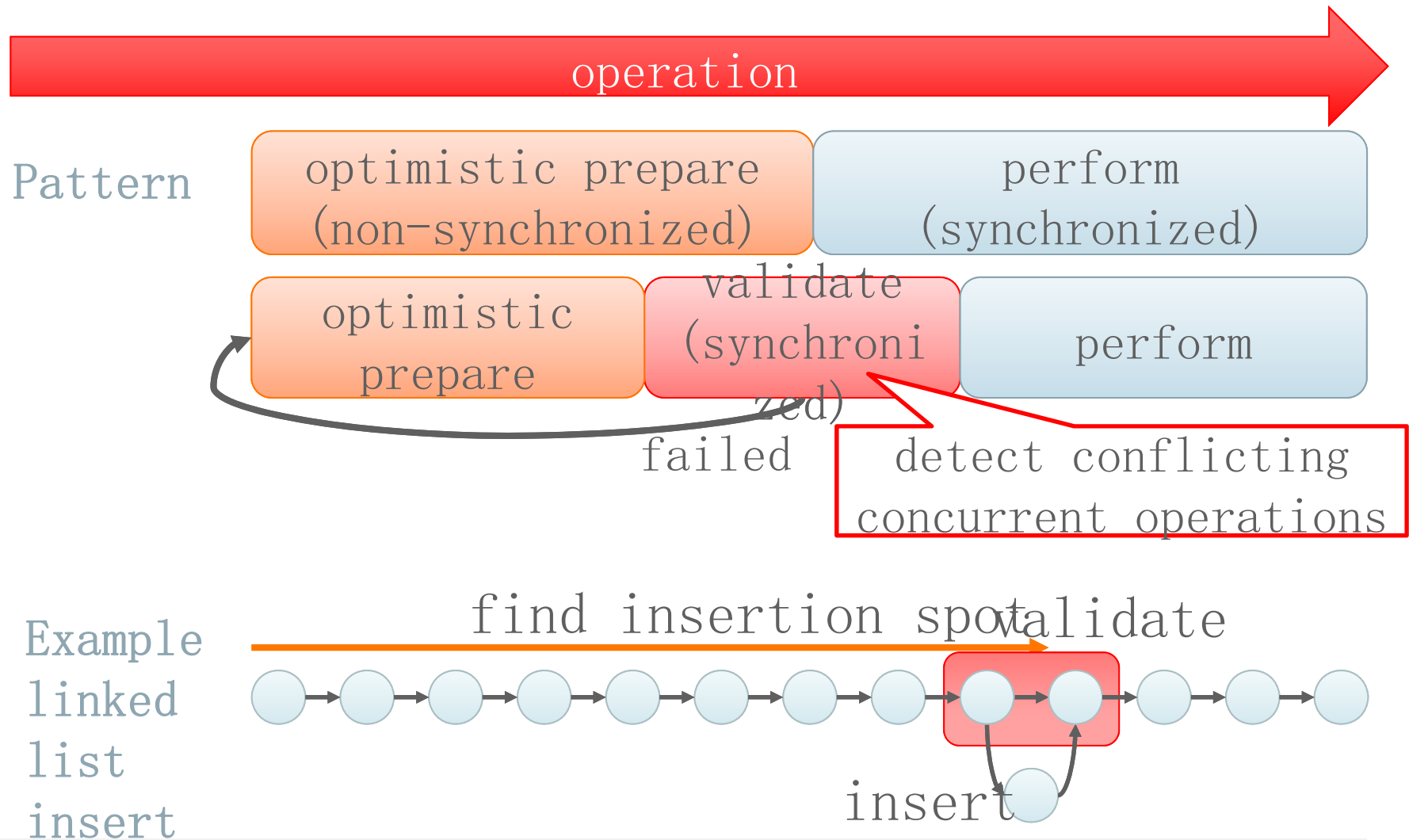# Concurrent Linked Lists – 5% Updates

Wait-free algorithm is slow ☹

# Optimistic Concurrency in Data Structures

**operation**

**Pattern**

| optimistic prepare (non-synchronized) | perform (synchronized) |

| optimistic prepare | validate (synchronized) | perform |

failed

detect conflicting concurrent operations

**Example linked list insert**

find insertion spot · validate

insert

Validation plays a key role in concurrent data structures

# Validation in Concurrent Data Structures

- **Lock-free**: atomic operations

  | optimistic prepare | validate & perform (atomic ops) |
  |---|---|

  failed

  – marking pointers, flags, helping, ⋯

- **Lock-based**: lock → validate

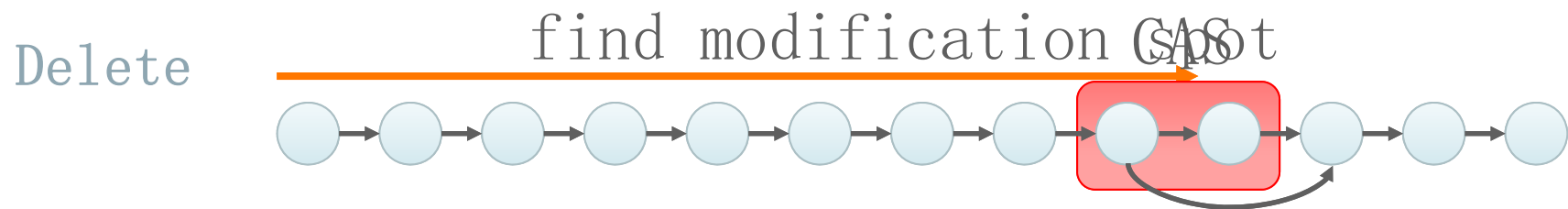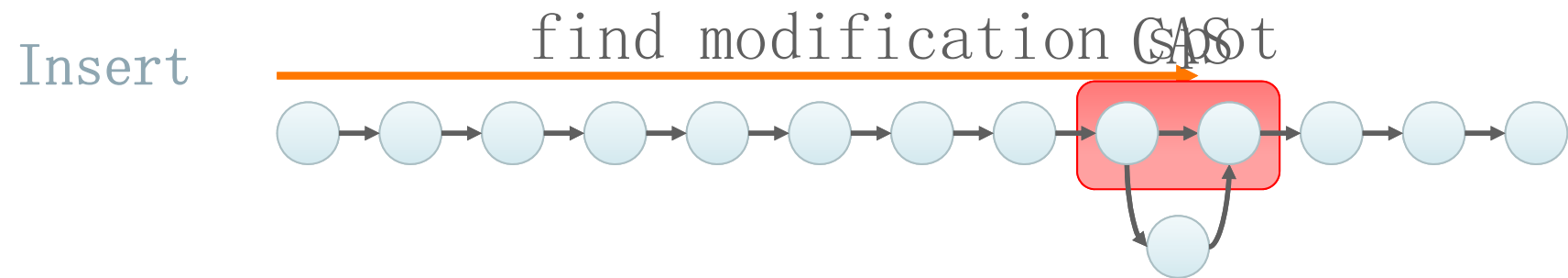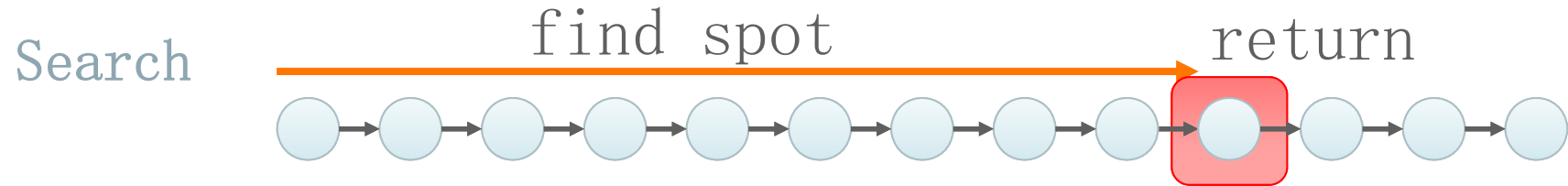  | optimistic prepare | lock | validate | perform | unlock |
  |---|---|---|---|---|

  unlock

  failed

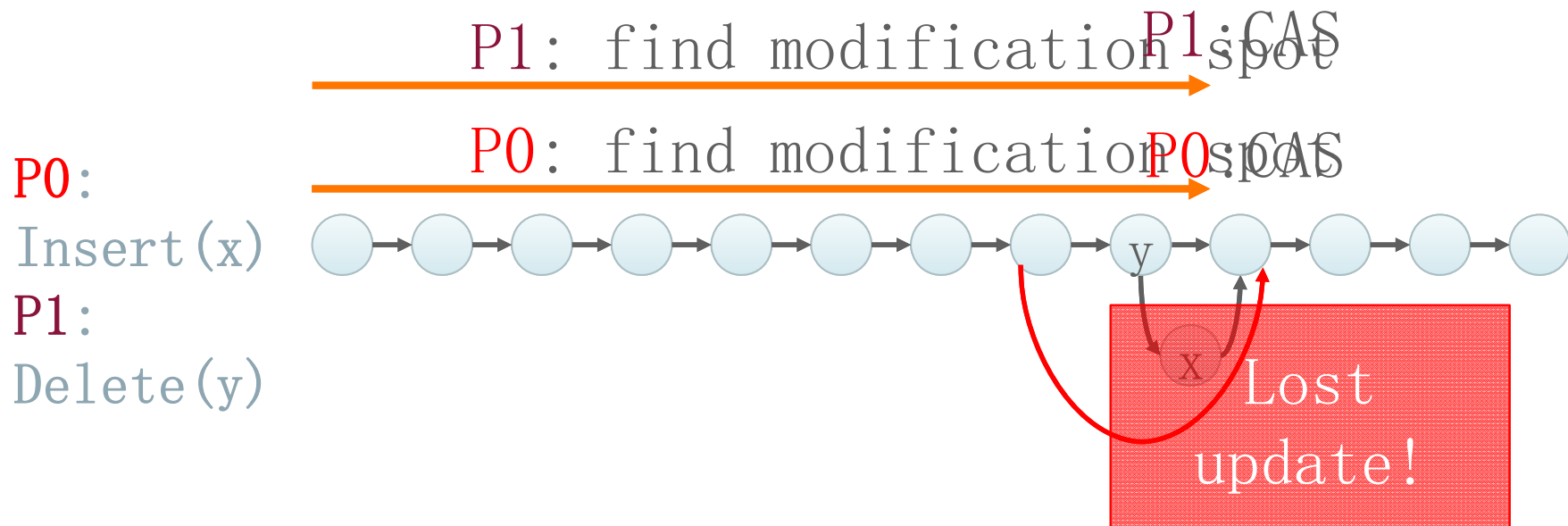  – flags, pointer reversal, parsing twice,

  Validation is what differentiates algorithms

# Let's design two concurrent linked lists:
# A lock-free and a lock-based

# Lock-free Sorted Linked List: Naïve

Search    find spot                return

Insert    find modification CAS spot

Delete    find modification CAS spot

Is this a correct (linearizable) linked list?

# Lock-free Sorted Linked List: Naïve – Incorrect



P1: find modification spot    P1:CAS

PO: find modification spot    P0:CAS

P0:
Insert(x)

P1:
Delete(y)

Lost update!

- What is the problem?
  - Insert involves one existing node;
  - Delete involves two existing nodes

How can we fix the problem?
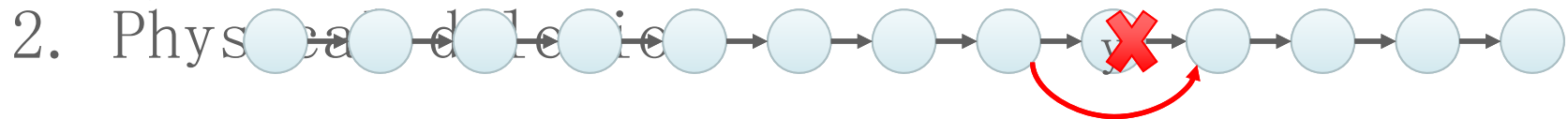
- **Idea**! To delete a node, make it unusable first···
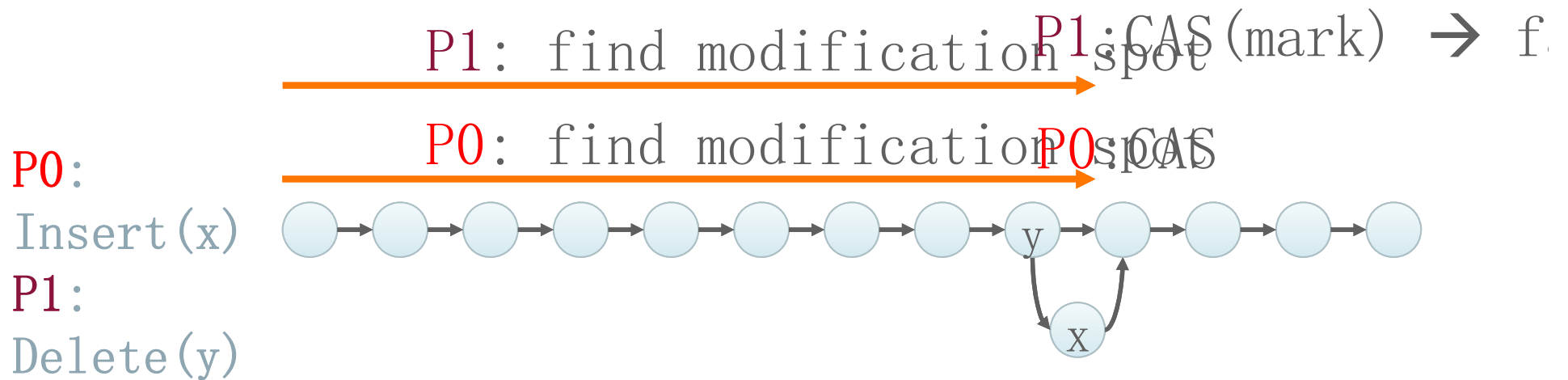
  - Mark it for **deletion** so that

    1. You fail marking if someone changes **next** pointer;

    2. An insertion fails if the predecessor node is marked.

→ In other words: delete in two steps

1. Mark for deletion; and then

2. Physical delete

Delete(y)

find modification spot

2. CAS(remove)

1. CAS(mark)

# 1. Failing Deletion (Marking)

P1: find modification spot    P1:CAS(mark) → f

P0: find modification spot    P0:CAS

PO:
Insert(x)

P1:
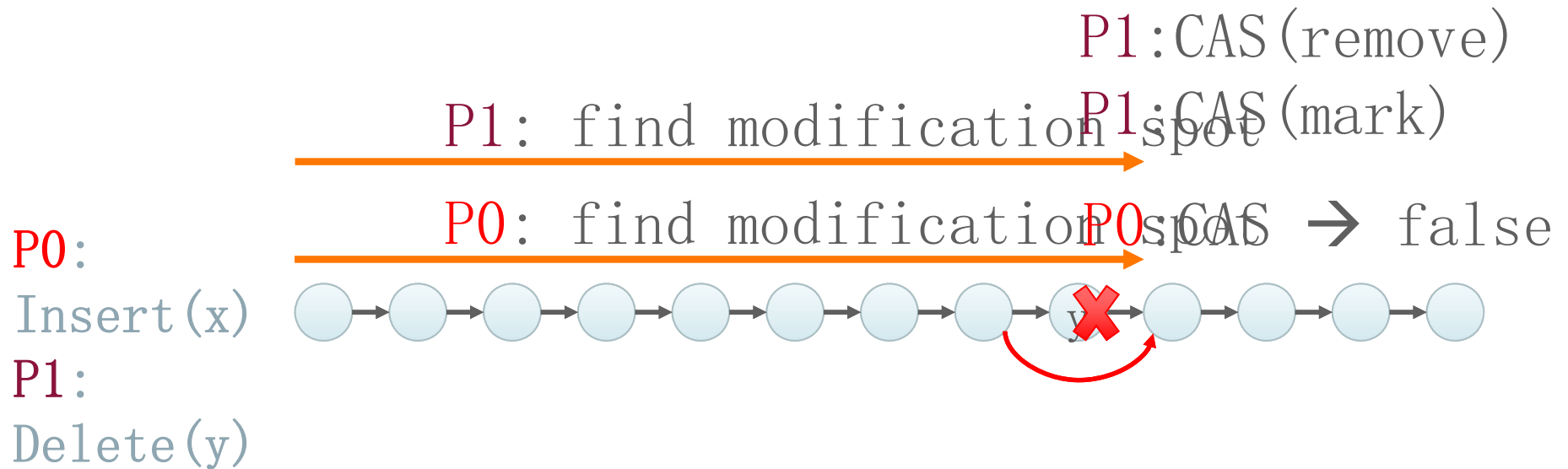Delete(y)



- Upon failure → restart the operation
  - Restarting is part of "all" state-of-the-art-data structures

# 1. Failing Insertion due to Marked Node

P1:CAS(remove)

P1:CAS(mark)

P1: find modification spot

P0: find modification spot  P0:CAS → false
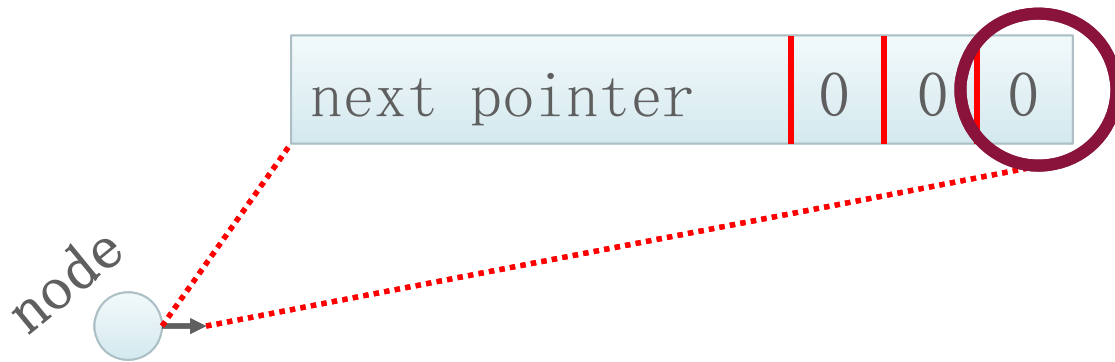
P0:
Insert(x)

P1:
Delete(y)

- Upon failure → restart the operation
  - Restarting is part of "all" state-of-the-art-data structures

How can we implement marking?

# Implementing Marking (C Style)

- Pointers in 64 bit architectures
  - Word aligned – 8 bit aligned!



```
boolean mark(node_t* n)
    uintptr_t unmarked = n->next & ~0x1L;
    uintptr_t marked   = n->next | 0x1L;
    return CAS(&n->next, unmarked, marked) == unmarked;
```
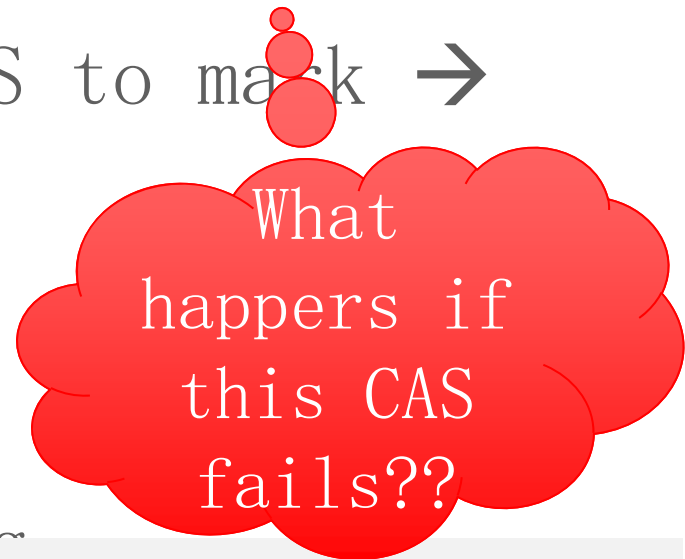
# Lock-free List: Putting Everything Together

- Traversal: traverse (requires unmarking nodes)

- Search: traverse

- Insert: traverse → CAS to insert

- Delete: traverse → CAS to mark → CAS to remove

- Garbage (marked) nodes
  - Cleanup while traversing
  - (*helping* in this course's terms)

What happers if this CAS fails??

A pragmatic implementation of lock-free linked lists

# What is not Perfect with the Lock-free List?

1. Garbage nodes
   - Increase path length; and
   - Increase complexity
     ```
     if (is_marked_node(n)) …
     ```
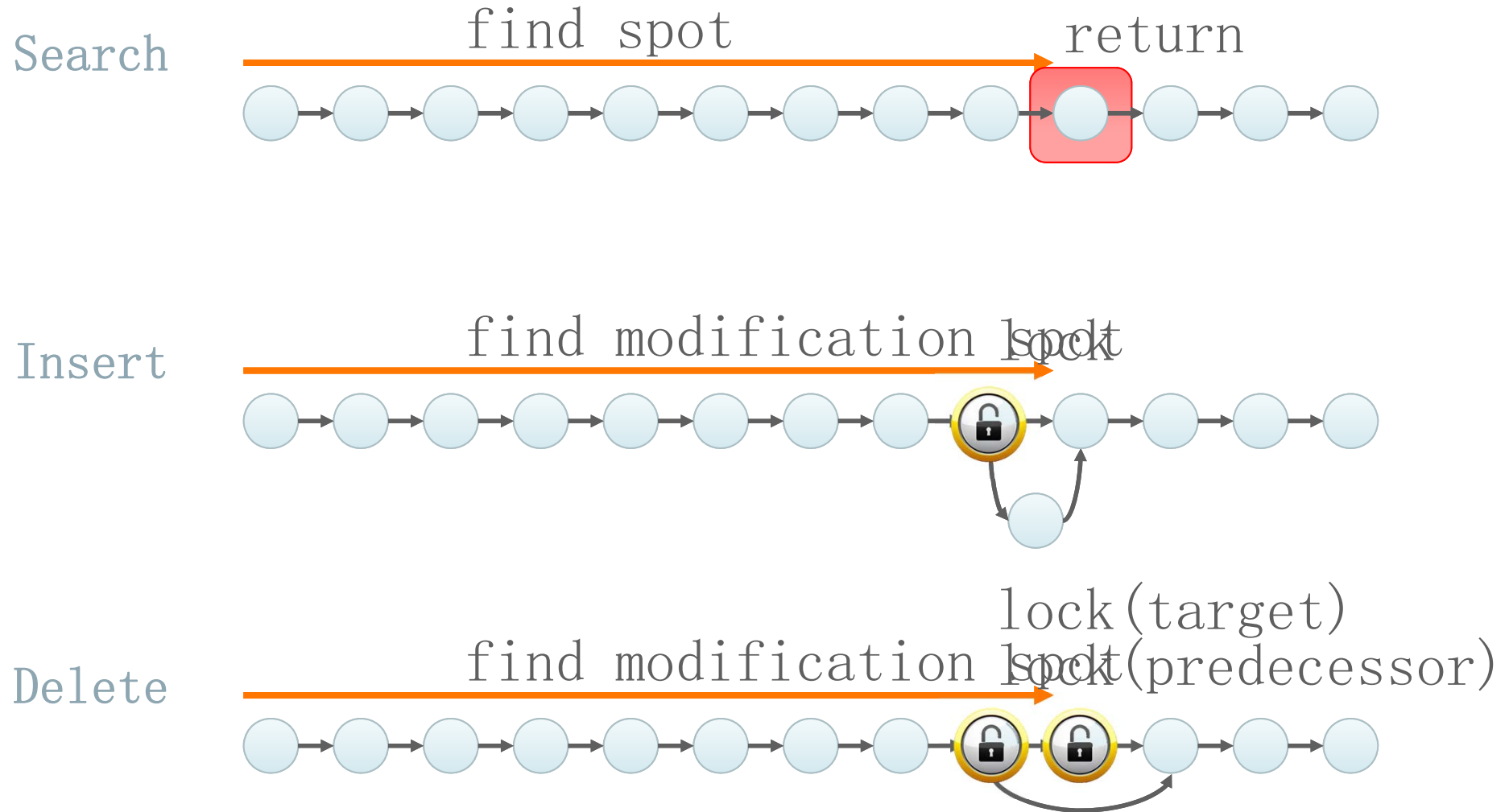
2. Unmarking every single pointer
   - Increase complexity
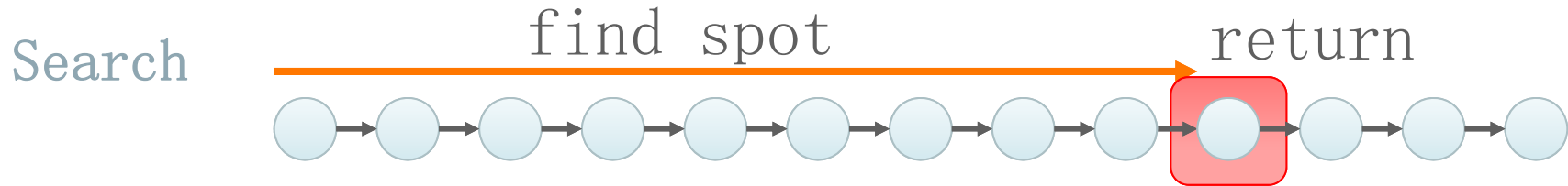     ```
     curr = get_unmark_ref(curr->next)
     ```

Can we simplify the design with locks?
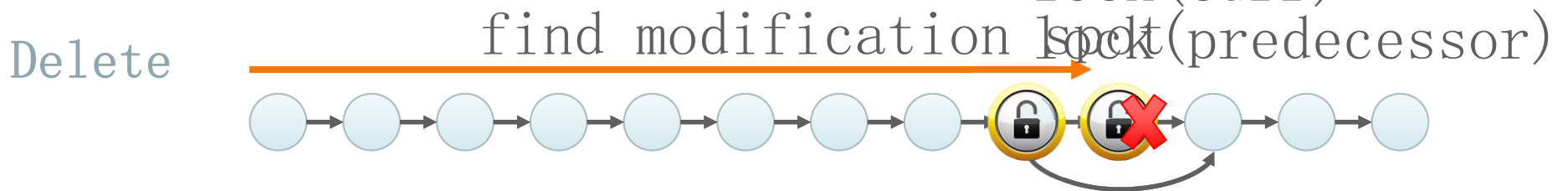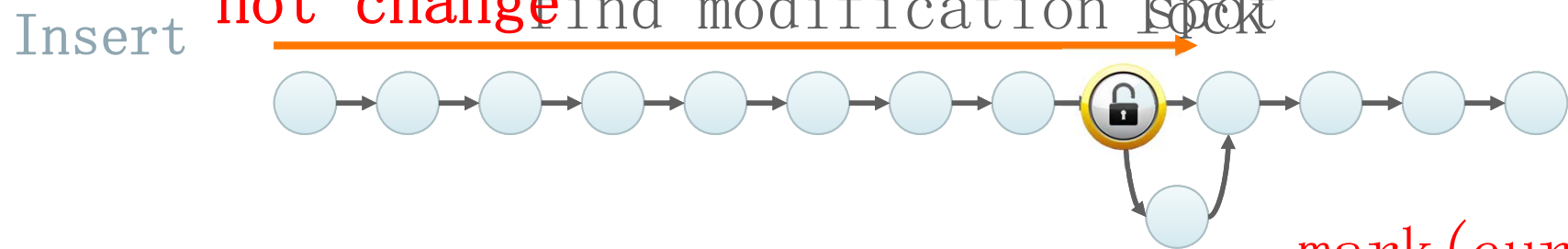
# Lock-based Sorted Linked List: Naïve

**Search**  find spot                return



**Insert**  find modification spotlock



lock(target)

**Delete**  find modification spotlock(predecessor)



Is this a correct (linearizable) linked list?

# Lock-based List: Validate After Locking

Search

find spot                    return

validate !pred->marked && pred->next did
not change find modification spot lock

Insert

mark(curr)
lock(curr)

find modification lock(predecessor)

Delete

!pred->marked && !curr->marked && pred->next did
change

# Concurrent Linked Lists – 0% updates

Just because the lock-based is not unmarking!

Throughput (Mop/s) vs # Cores

- - - lock-free — lock-based

(Lesson$_2$) Sequential complexity matters → Simplicity
☺

# Optimistic Concurrency Control: Summary

- **Lock-free**: atomic operations

  | optimistic prepare | validate & perform (atomic ops) |
  |---|---|

  failed

  – marking pointers, flags, helping, …

- **Lock-based**: lock → validate

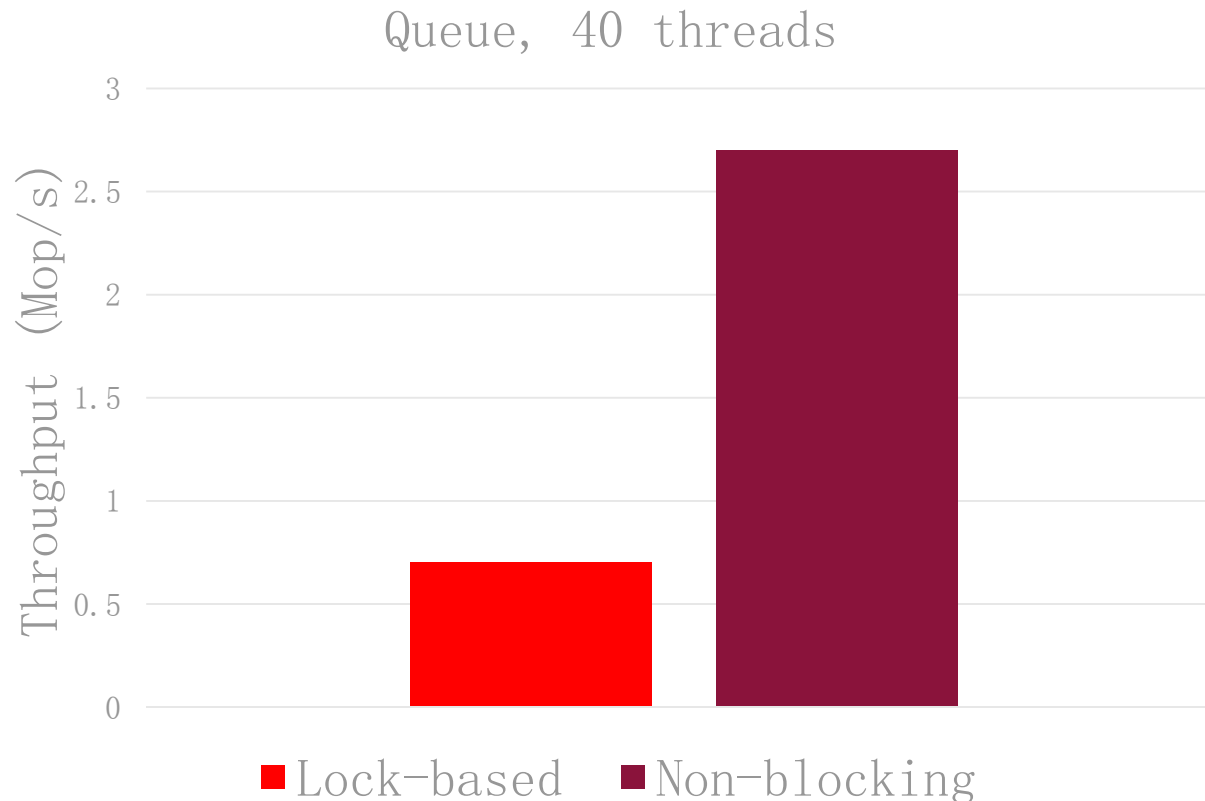  | optimistic prepare | lock | validate | perform | unlock |
  |---|---|---|---|---|

  unlock

  failed

  – flags, pointer reversal, parsing twice, …
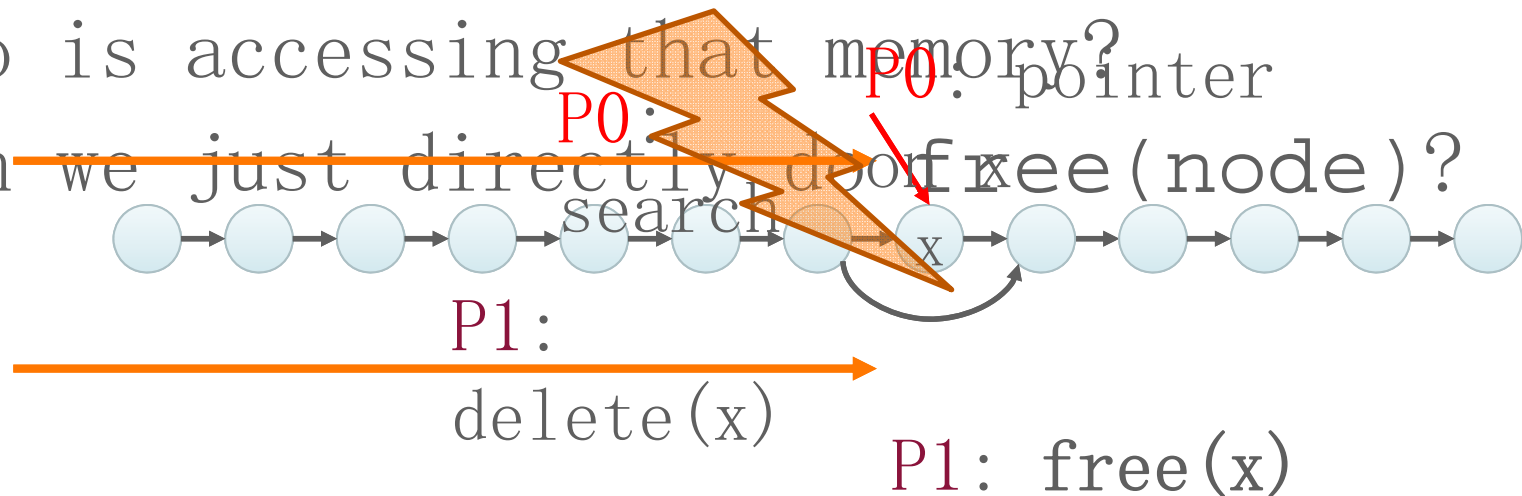
# Word of caution: lock-based algorithms

- Search data structures   ☺
- Queues, stacks, counters, ...   ☹

Queue, 40 threads

# Memory Reclamation: OCC's Side Effect

- Delete a node → free and reuse this memory

- Subset of the garbage collection problem

- Who is accessing that memory?

- Can we just directly do free(node)?

P0: search

P0: pointer

P1: delete(x)

P1: free(x)

x

We cannot directly free the memory! Need memory reclamation

# Memory Reclamation Schemes

1. ## Reference counting
   - Count how many references exist on a node

2. ## Hazard pointers
   - Tell to others what exactly you are reading

3. ## Quiescent states
   - Wait until it is certain than no one holds references

4. ## Read-Copy Update (RCU)
   - Quiescent states – The extreme approach

# 1. Reference Counting

- Pointer + Counter

rc_pointer

counter | pointer

- Dereference:
  **rc_dereference**(rc_pointer* rcp)
      atomic_increment(&rcp->counter);
      return *pointer;

- "Release":
  **rc_release**(rc_pointer* rcp)
      atomic_decrement(&rcp->counter);

- Free: iff counter = 0

$(Lesson_3)$ Readers cannot write on the shared nodes

Bad bad bad idea: Readers write on shared nodes!
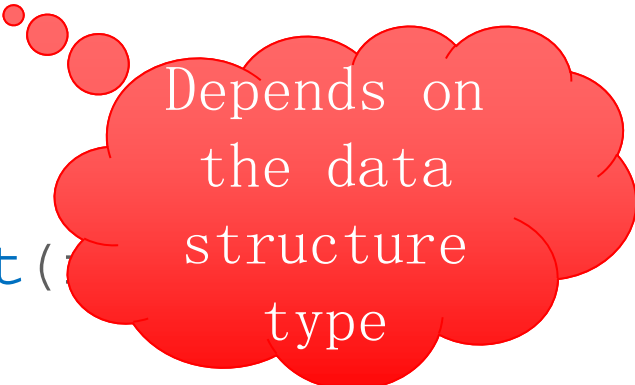
# 2. Hazard pointers (1/2)

- Reference counter → property of the node
- Hazard pointer → property of the **hazard_pointer address** thread
  - A Multi-Reader Single-Writer (MRSW) register
- Protect:
  **hp_protect**(node* n)
      hazard_pointer* hp = hp_get(
      hp->address = n;

  > Depends on the data structure type

- Release:
  **hp_release**(hazard_pointer* hp)

# 2. Hazard pointers (2/2)

- Free memory **x**

  1. Collect all hazard pointers
  2. Check if **x** is accessed by any thread hazard_pointe
     1. If yes, buffer the free for later <span style="background:red;color:white">addres</span>
     2. If not, free the memory

- Buffering the free is implementation specific

- + lock-free

O(data structure size) hazard pointers hp_protect

# 3. Quiescent States

- Keep the memory until it is certain it is not accessed

- Can be implemented in various ways

- Example implementation
  **search / insert / delete**
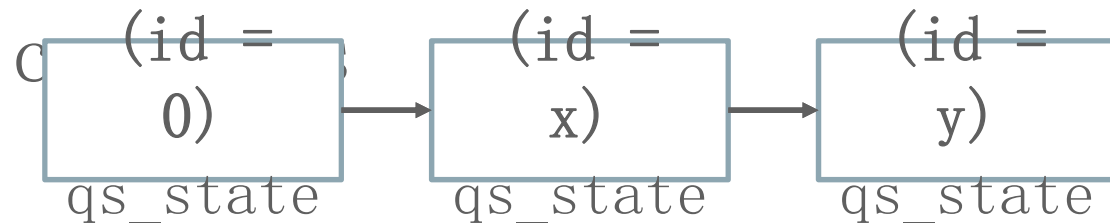
  qs_unsafe();     I'm accessing shared data

  …

  qs_safe();       I'm not accessing shared data

  The data written in qs_[un]safe must be local-mostly

- List of "thread-local" (mostly)

```
c   (id =        (id =        (id =
       0)           x)           y)
    qs_state     qs_state     qs_state
```
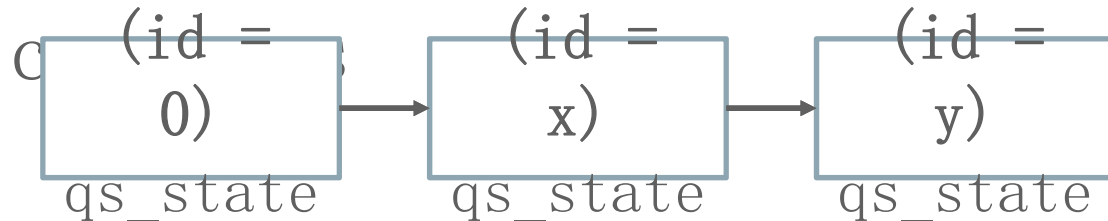
- qs_state (initialized to 0)
  - even : in safe mode (not accessing shared data)
  - odd : in unsafe mode
- qs_safe / qs_unsafe

How do we free memory?

# 3. Quiescent States: Freeing memory

- List of "thread-local" (mostly)



| (id = 0) | (id = x) | (id = y) |
| --- | --- | --- |
| qs_state | qs_state | qs_state |

- Upon `qs_free`: Timestamp memory (`vector_ts`)
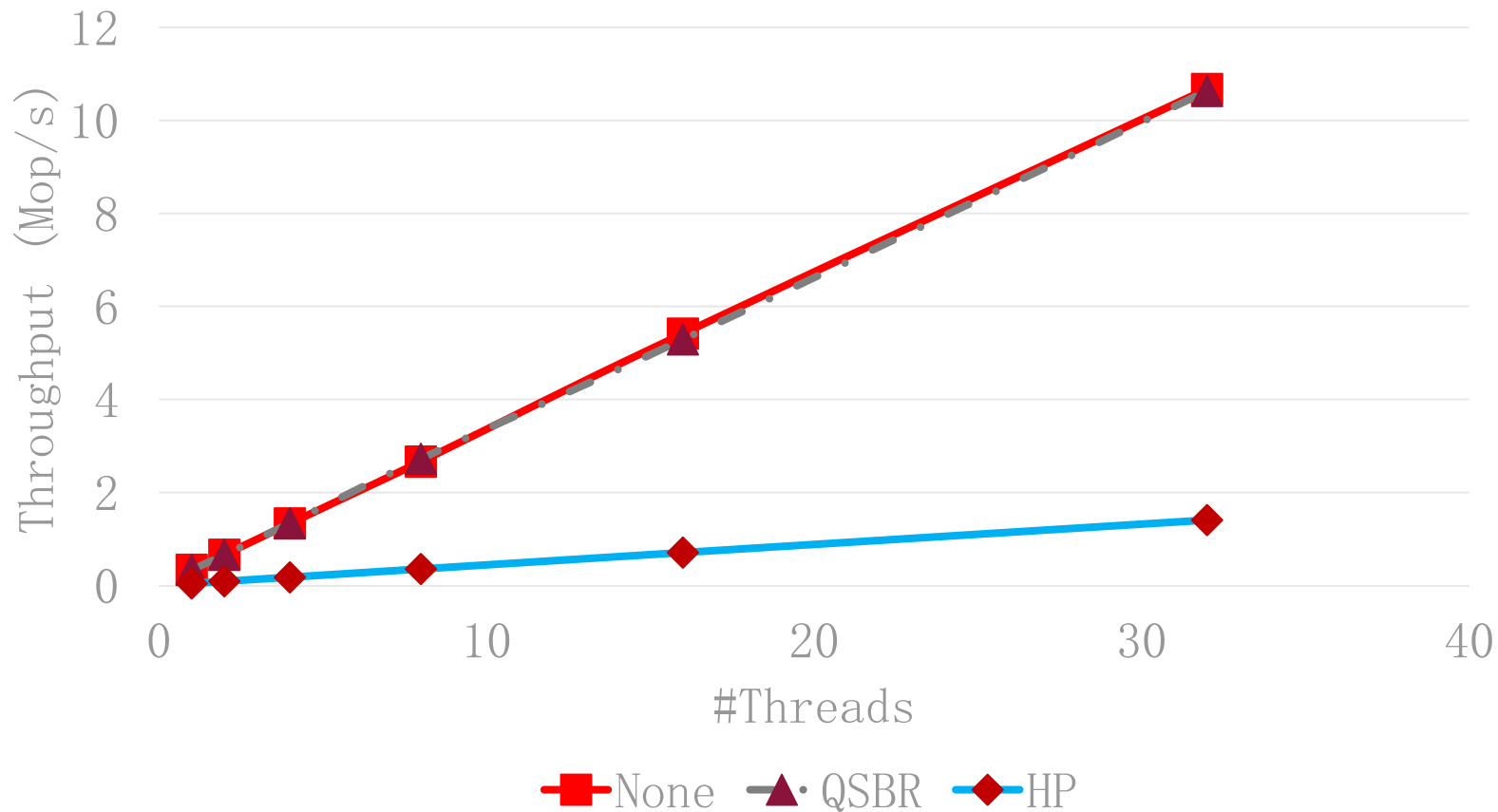  - Can do this for batches of frees

```
for t in thread_ids
    if (vts_mem[t] is odd &&
        vts_now[t] = vts_mem[t])
        return false;
return true;
```

- Safe to reuse the memory
  vector_ts   >> vector_ts

How do the schemes we have seen perform?

# Hazard Pointers vs. Quiescent States

Quiescent-state reclamation is as fast as it gets

# 4. Read-Copy Update (RCU)

- Quiescent states at their extreme
  - Deletions <span style="color:red">wait all readers</span> to reach a safe state
- Introduced in the Linux kernel in ~2002
  - More than 10000 uses in the kernel!
- (Example) Interface
  - `rcu_read_lock` (= `qs_unsafe`)
  - `rcu_read_unlock` (= `qs_safe`)
  - `synchronize_rcu` → wait all readers

# 4. Using RCU

- **Search / Traverse**
  `rcu_read_lock()`
  …
  `rcu_read_unlock()`

- **Delete**
  … physical deletion of **x**
  `synchronize_rcu()`
  `free(`**x**`)`

- **+** simple
- **+** read-only workloads
- **−** bad for writes

# Memory Reclamation: Summary

- How and when to reuse freed memory
- Many techniques, no silver bullet
  1. Reference counting
  2. Hazard pointers
  3. Quiescent states
  4. Read-Copy Update (RCU)

# Summary

- Concurrent data structures are very important

- Optimistic concurrency necessary for scalability
  - Only recently a lot of active work for CDSs

- Memory reclamation is
  - Inherent to optimistic concurrency;
  - A difficult problem;
  - A potential performance/scalability bottleneck