# Programmazione concorrente

Laurea Magistrale in Ingegneria Informatica
Università Tor Vergata
Docente: Romolo Marotta
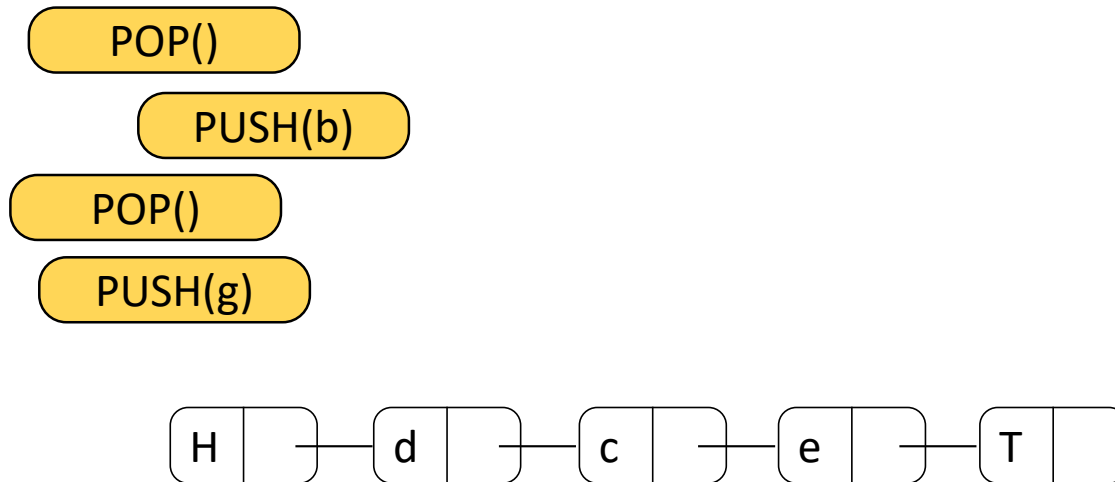
# Concurrent data structures

1. Stack
2. Set
3. Priority queues
4. FIFO queues
5. MRSW Registers
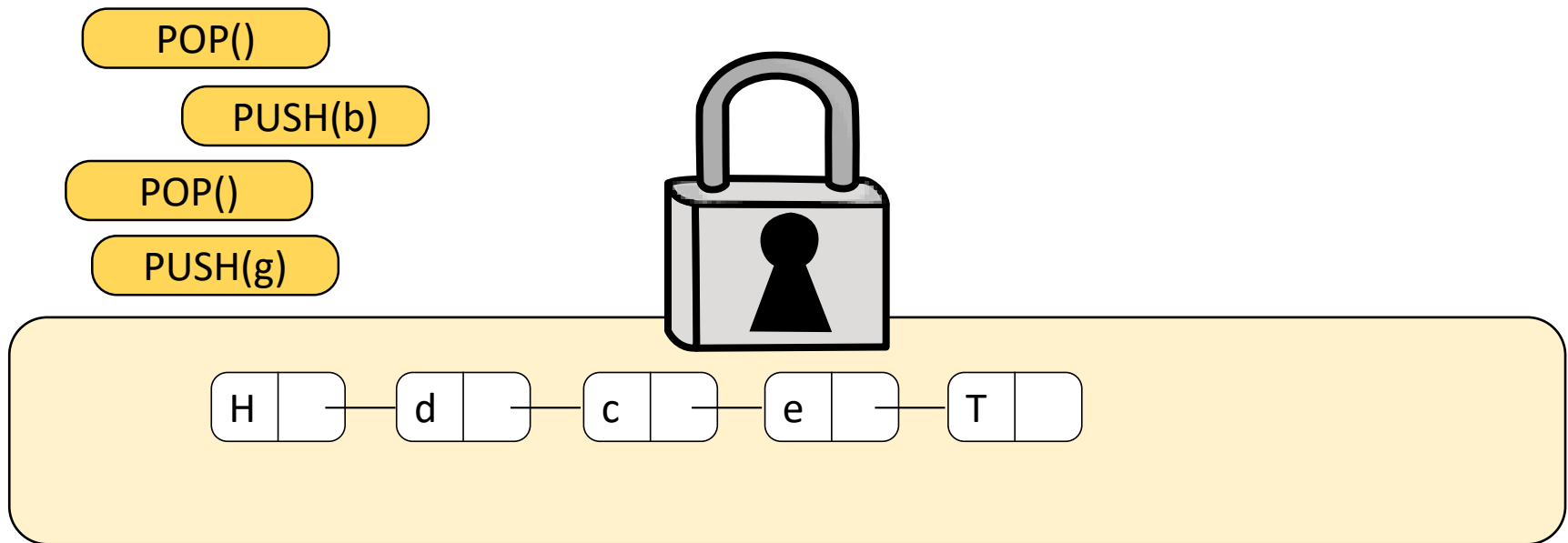
# Concurrent
# Data Structures:
# Stacks

# Stack implementation

- Stack methods:
  - push(v)
  - pop()
- Implemented as a linked list

# Concurrent stack implementations

- Resort to a global lock

# Read-Modify-Write

- RMW instructions allow to read memory and modify its content in an apparently instantaneous fashion.

```
1.RMW(MRegister *r, Function f){
2.   atomic{
3.     old = r;
4.     *r = f(r);
5.     return old;
6.   }
7.}
```

- Even conventional atomic Load and Store can be seen as RMW operations

# Compare-And-Swap

- Compare-and-Swap (CAS) is an atomic instruction used in multithreading to achieve synchronization
  - It compares the contents of a memory area with a supplied value
  - If and only if they are the same
  - The contents of the memory area are updated with the new provided value
- Atomicity guarantees that the new value is computed based on up-to-date information
- If, in the meanwhile, the value has been updated by another thread, the update fails
- This instruction has been introduced in 1970 in the IBM 370 trying to limit as much as possible the use of spinlocks

# Compare-And-Swap

- RMW instructions allow to read memory and modify its content in an apparently instantaneous fashion.
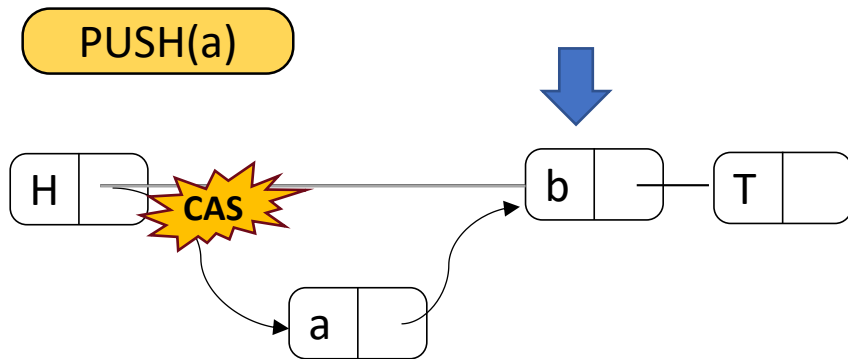
```
1. CAS(Mregister *r, Value old_value, Value new_value f){
2.    atomic{
3.        Value res = *r;
4.        if(*r == old_value) *r = new_value;
5.        return res;
6.    }
7. }
```

- CAS is implemented by x86 architectures (see CMPXCHG)
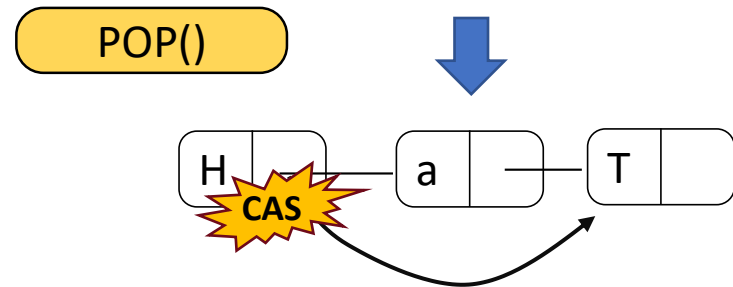- gcc offers the __sync_val_compare_and_swap builtin

# Attempt 1

Push:

1. Get head next
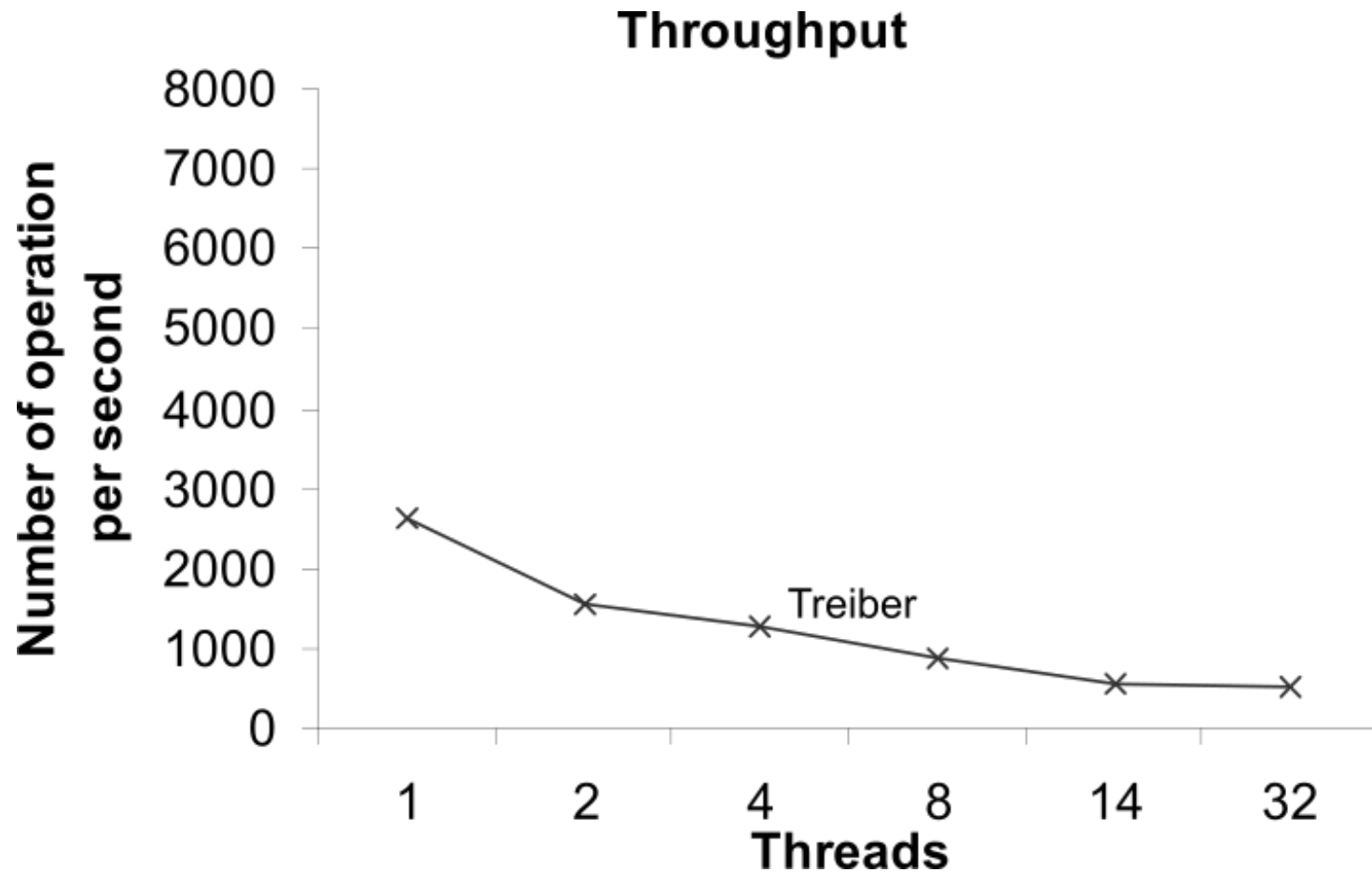2. Insert the new item with a CAS
3. If CAS fails, restart

Delete:

1. Get head next
2. Disconnect the item with a CAS
3. If CAS fails, restart

PUSH(a)

H — CAS — b — T

a

POP()

H — CAS — a — T
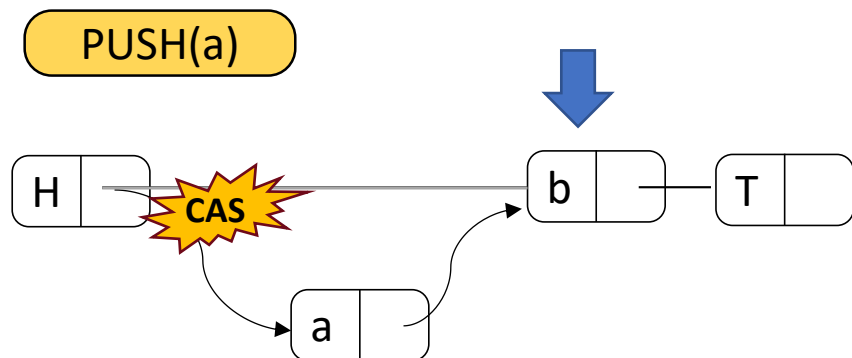
- Is it scalable?

# Non-blocking stack – Attempt 2 [Treiber+BO]

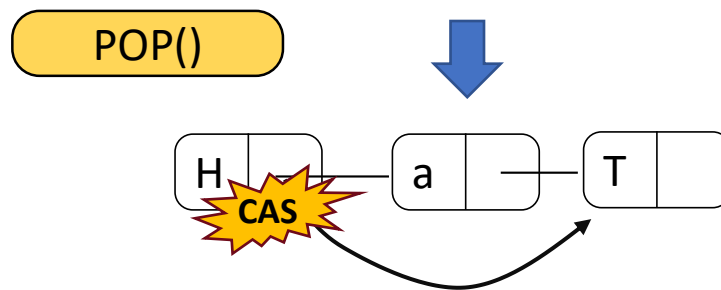# Non-blocking stack – Attempt 2 [Treiber+BO]

Push:

1. Get head next

2. Insert the new item with a CAS

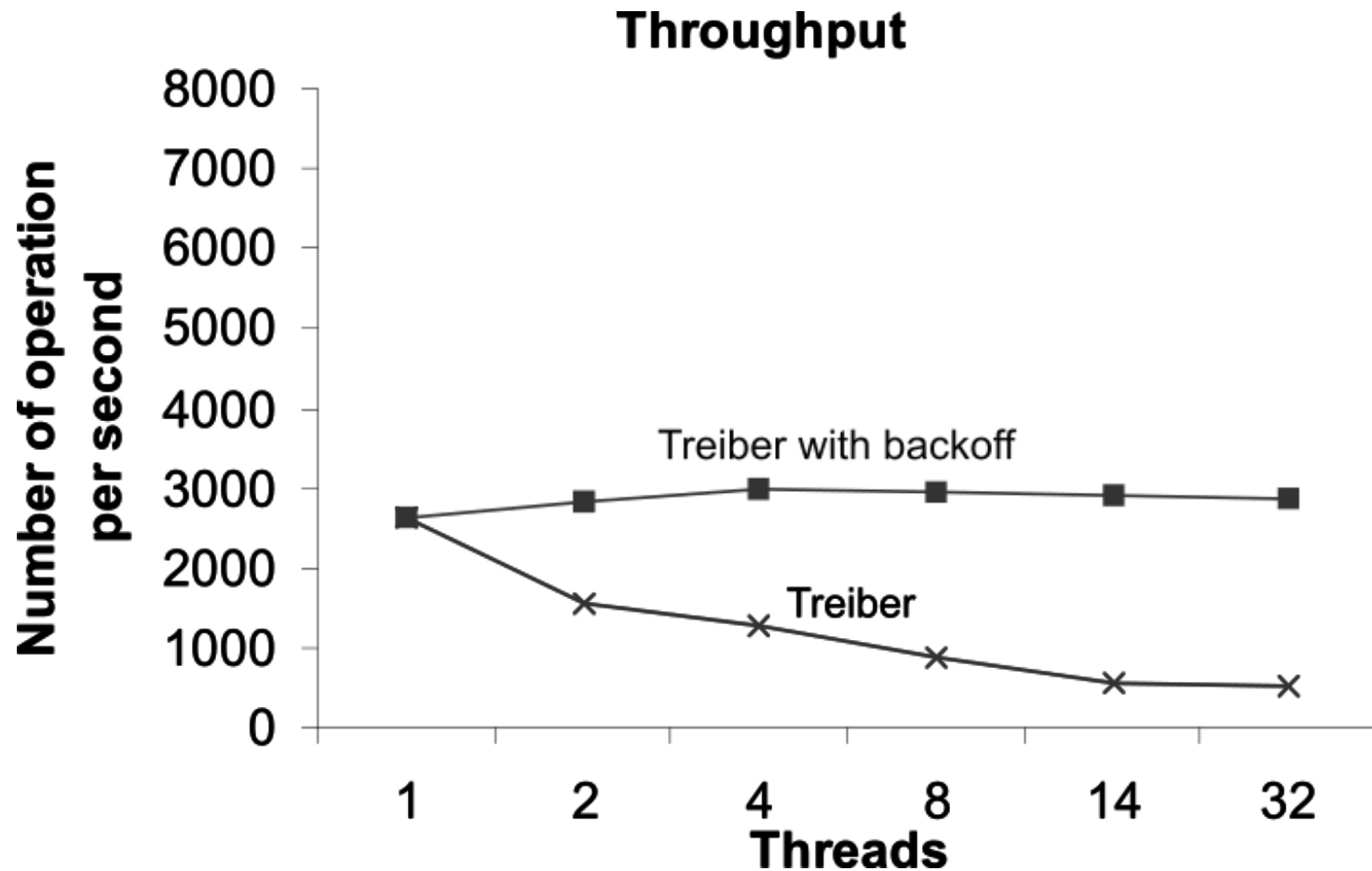3. If CAS fails, ~~restart~~ backoff and restart

Delete:

1. Get head next

2. Disconnect the item with a CAS

3. If CAS fails, ~~restart~~ backoff and restart

PUSH(a)

H | → b | — T

CAS

a |

POP()

H | — a | — T

CAS

- Is it scalable?

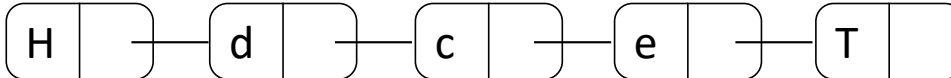# Non-blocking stack – Attempt 2 [Treiber+BO]



**Throughput**

# Concurrent stack implementations

- Resort to a global lock
  - Do not scale
- Resort to a naïve non-blocking approach
  - Do not scale
- Resort to a naïve non-blocking approach + Back off
  - Do not scale, but conflict resilient
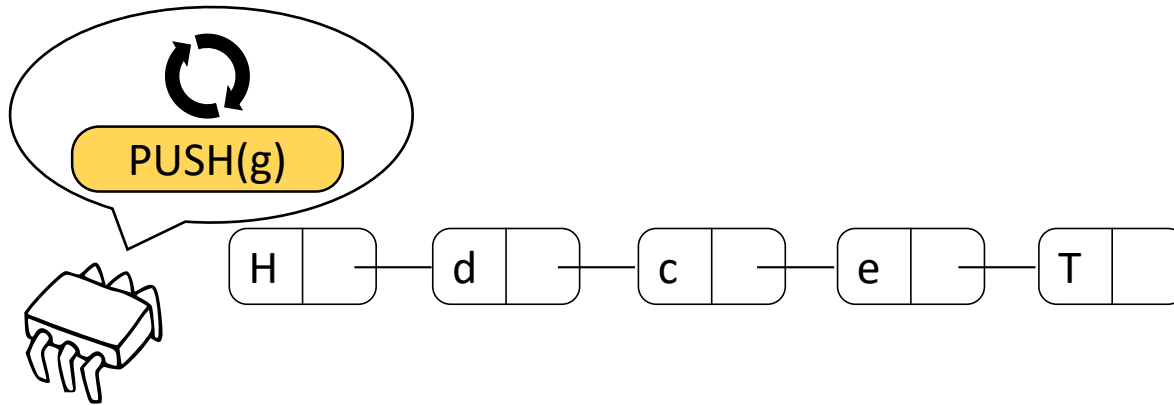- How achieve scalability? Make back-off times useful
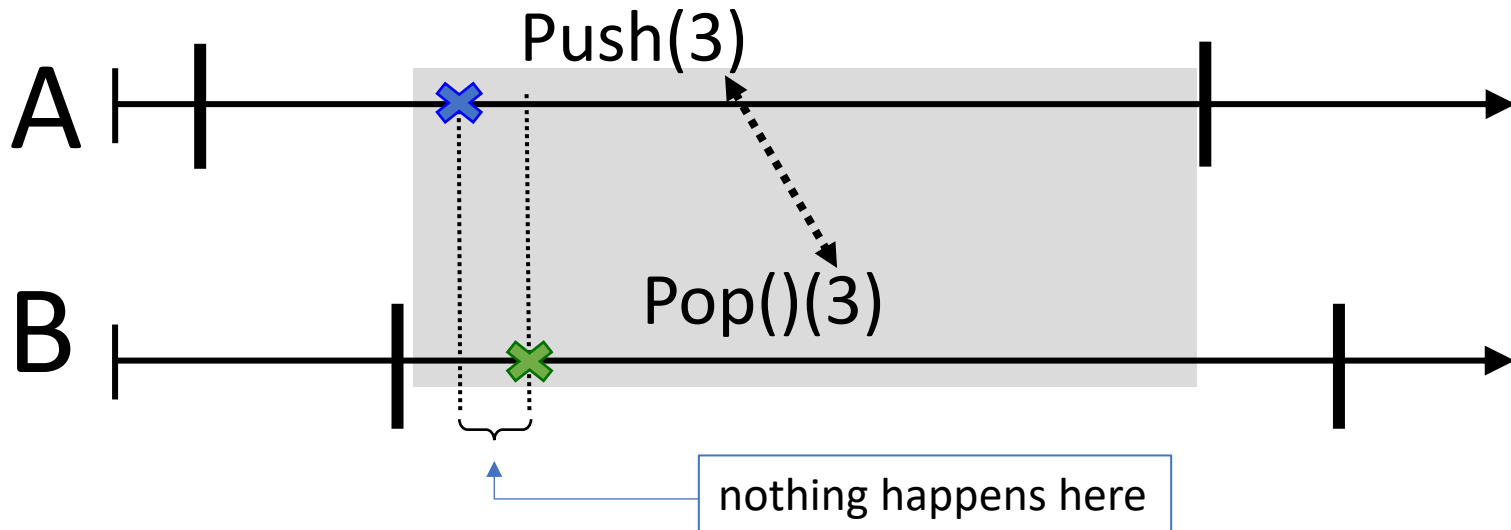
POP()

PUSH(g)

H — d — c — e — T

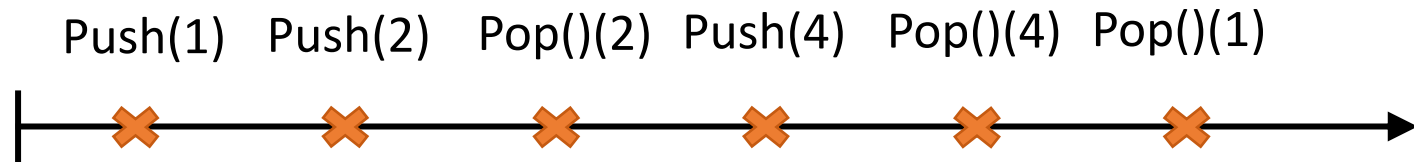# Non-blocking stack – Attempt 3

- How to take advantage of back-off times?

# Observation

• Concurrent matching push/pop pairs are always linearizable

Push(3)

A

Pop()(3)

B

nothing happens here

• A push A and a pop B are:
  ◦ concurrent to each other
  ◦ B returns the item inserted by A

⇒ we can always take two points such that:
  ◦ A is the last one to insert an item before A linearizes
  ◦ B appears to extract the last item inserted (by A)

# Observation

- Concurrent matching push/pop pairs are always linearizable

Push(3)

A

Pop()(3)

B

Push(1)  Push(2)  Pop()(2)  Push(4)  Pop()(4)  Pop()(1)

# Observation

- Concurrent matching push/pop pairs are always linearizable

Push(3) Pop()(3)

Push(1)  Push(2)  Pop()(2)  Push(4)  Pop()(4)  Pop()(1)

# Non-blocking stack – Attempt 3
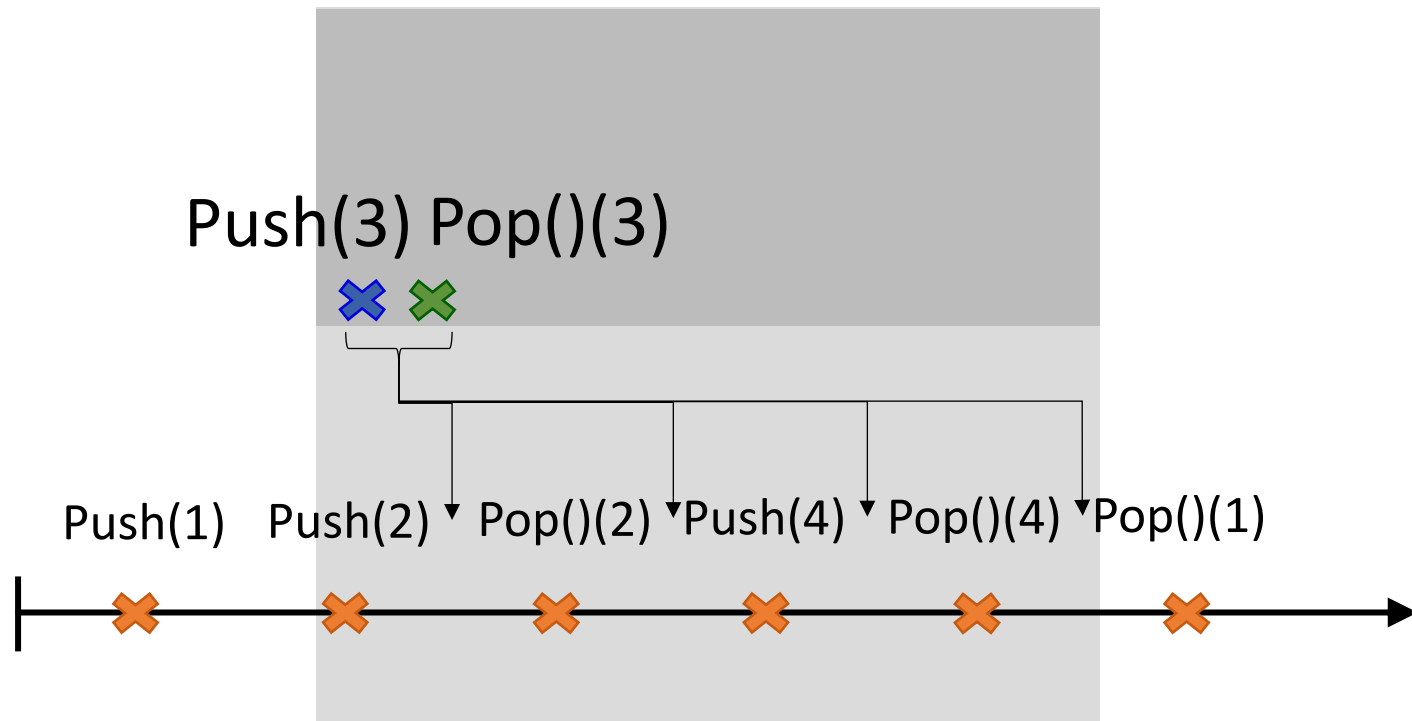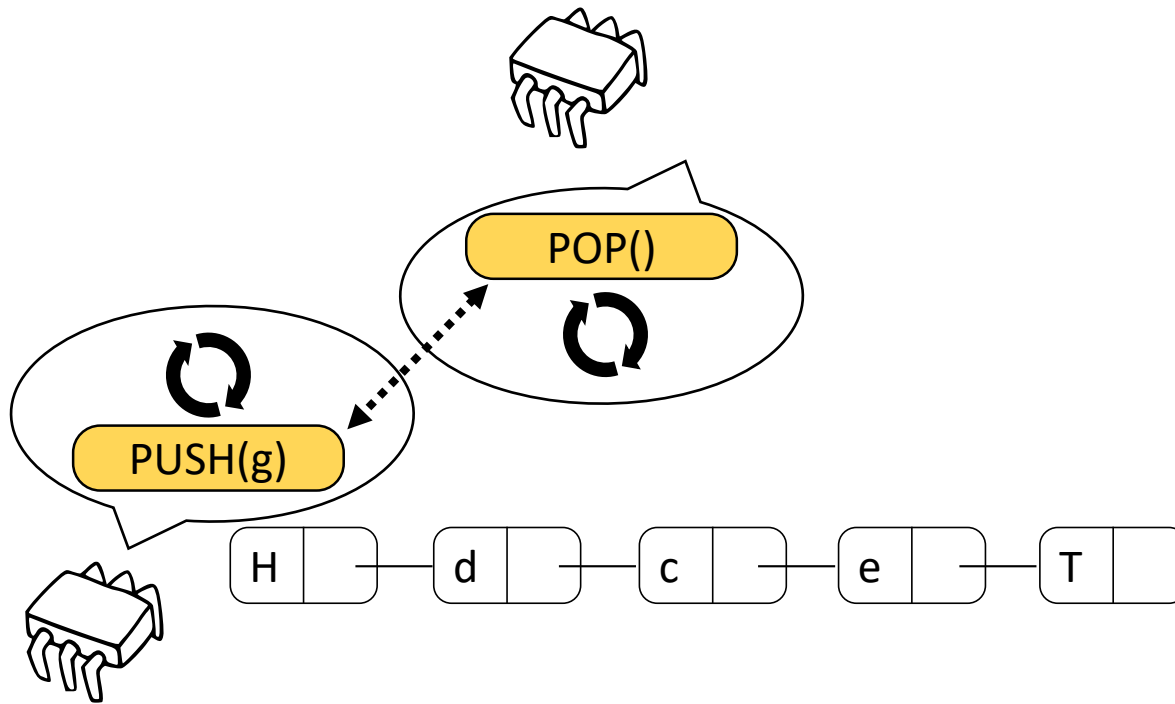
- How to take advantage of back-off times?
- Hope that an opposite operation arrives while waiting
- Match the two without interacting with the stack
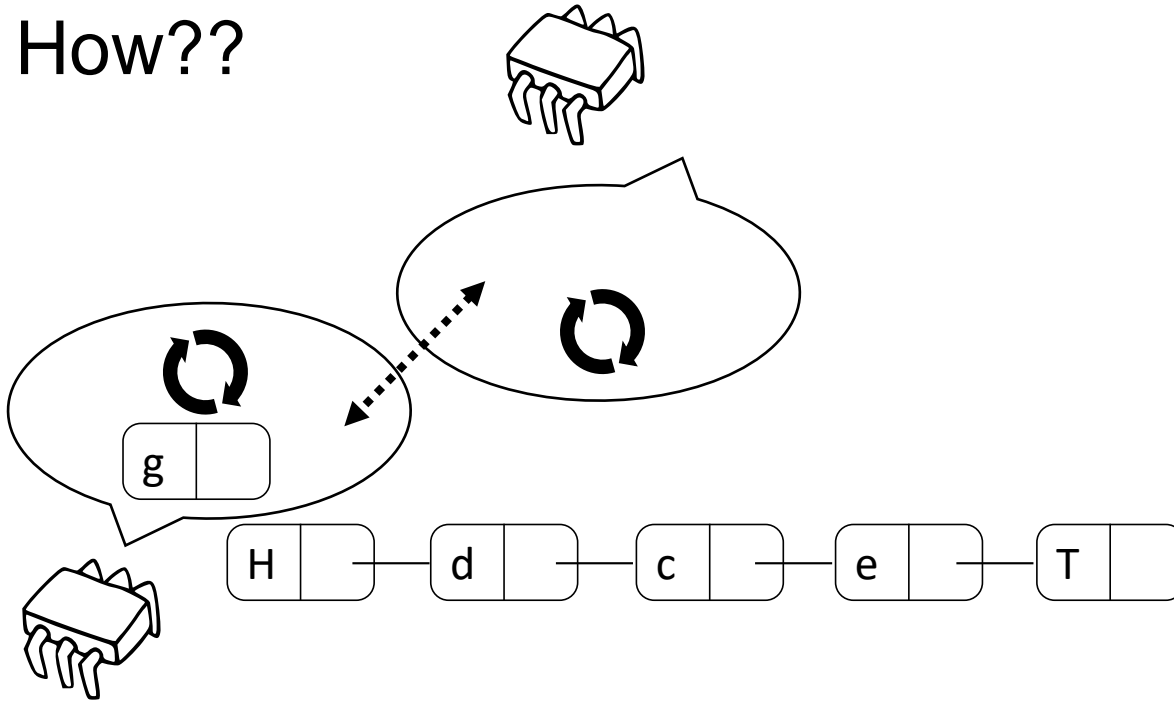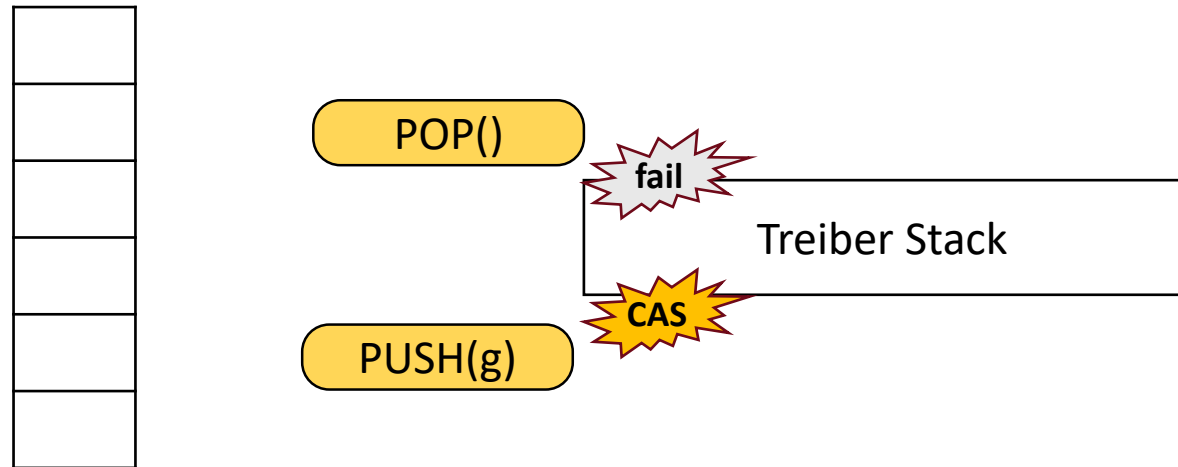
# Non-blocking stack – Attempt 3

- How to take advantage of back-off times?
- Hope that an opposite operation arrives while waiting
- Match the two without interacting with the stack
- How??

# Non-blocking stack – Elimination stack

- Pair the Treiber stack with an array

- Algorithm:
  1. Update the original stack via CAS
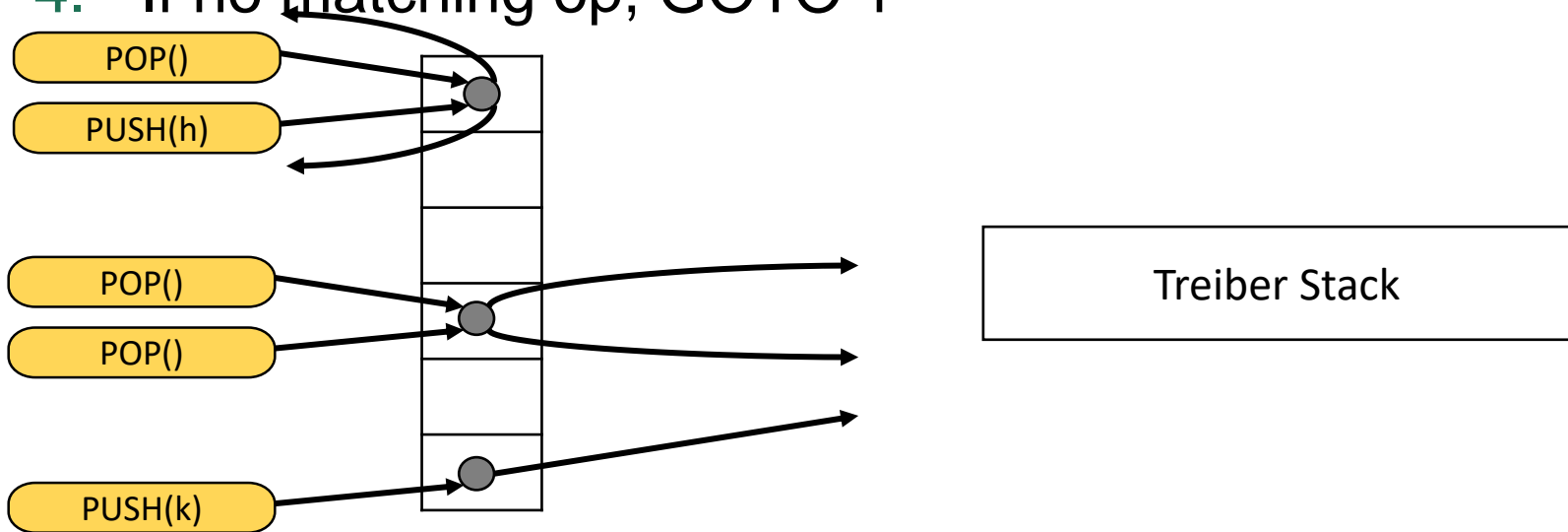  2. If CAS fails, publish the operation in a random cell of the array
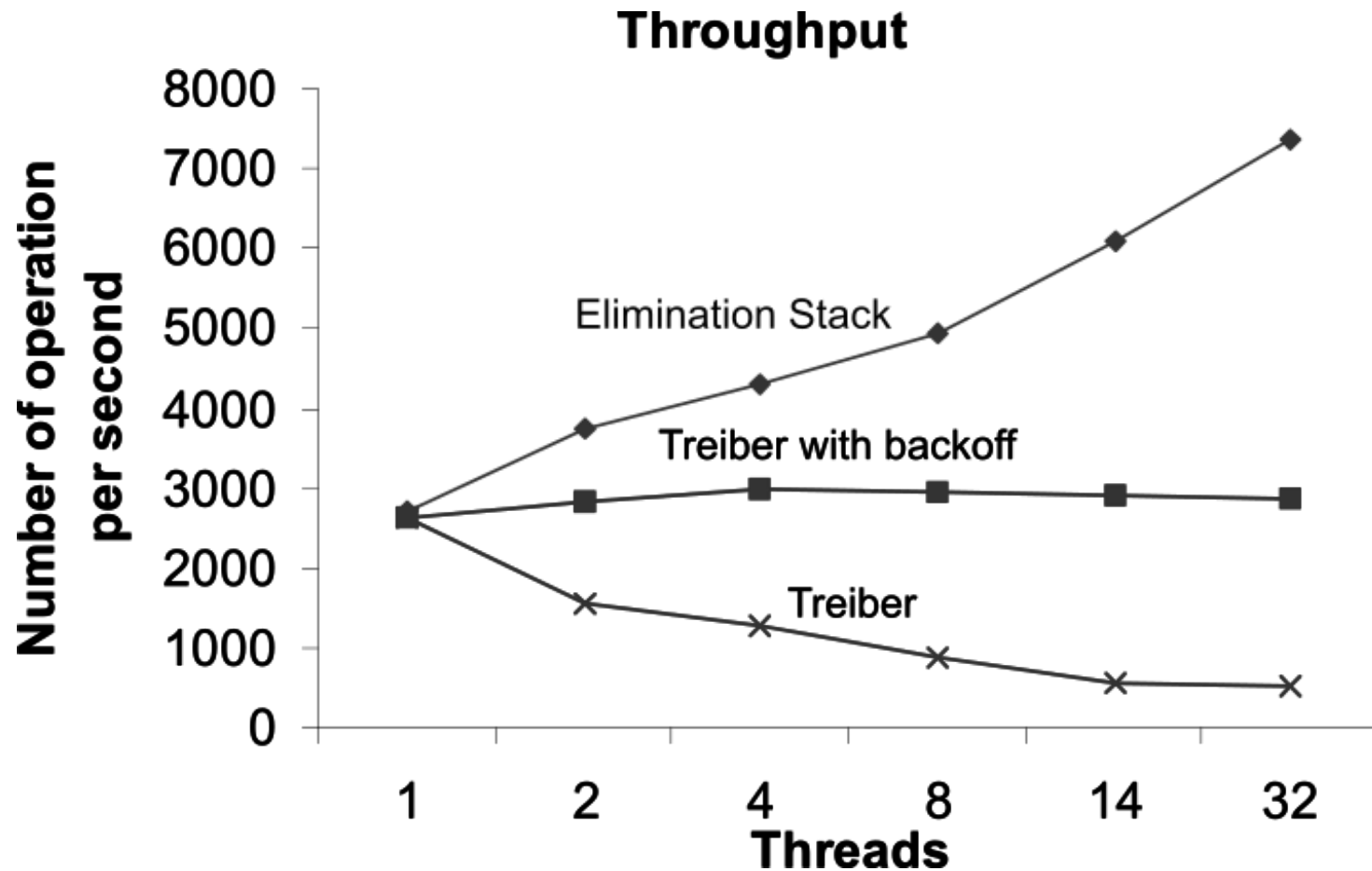
# Non-blocking stack – Elimination stack

- Pair the Treiber stack with an array
- Algorithm:
  1. Update the original stack via CAS
  2. If CAS fails, publish the operation in a random cell of the array
  3. Wait for a matching operation
  4. If no matching op, GOTO 1



POP()
PUSH(h)
POP()
POP()
PUSH(k)

Treiber Stack

# Non-blocking stack – Attempt 3



Throughput

# Concurrent
# Data Structures:
# Sets

# Set implementations

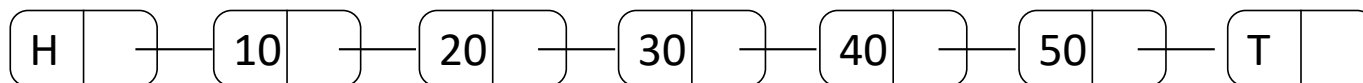- Set methods:
  - `insert(k)`
  - `delete(k)`
  - ~~`find(k)`~~

- Implemented as an ordered linked list

INSERT(35)

INSERT(25)

DELETE(40)

INSERT(55)

| H | | 10 | | 20 | | 30 | | 40 | | 50 | | T | |

# Insert algorithm

INSERT(55)

H — 10 — 20 — 30 — 40 — 50 — T

# Insert algorithm

# Insert algorithm

# Insert algorithm

left    right

H — 10 — 20 — 30 — 40 — 50 — 55 — T

# Delete algorithm

DELETE(40)

H → 10 → 20 → 30 → 40 → 50 → T

# Delete algorithm

left          right

H | 10 | 20 | 30 | 40 | 50 | T
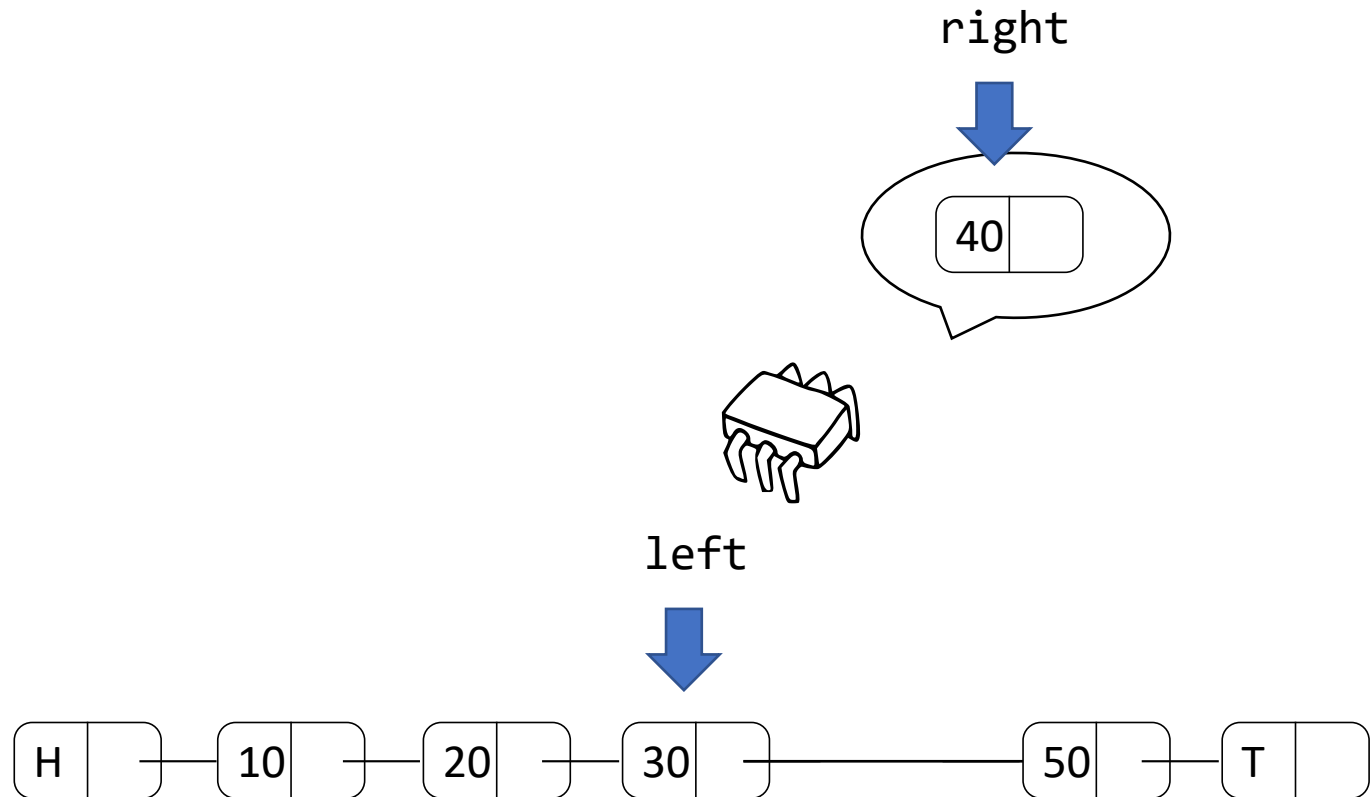
# Delete algorithm

# Sequential set implementation

```
1. bool do_operation(int k, int op_type){
2.    bool res = true;
3.    node *l,*r;
4.
5.    l = search(k, &r);
6.    switch(op_type){
7.      case(INSERT):
8.        if(r->key == k)
9.          res = false;
10.       else
11.         l->next = new node(k,r);
12.       break;
13.     case(DELETE):
14.       if(r->key == k)
15.         l->next = r->next;
16.       else
17.         res = false;
18.       break;
19.    }
20.
21.
22.    return res;
23.}
```

```
1. node* search(int k, node **r){
2.    node *l, *r_next;
3.    l = set->head;
4.
5.    *r = l->next;
6.
7.    r_next = (*r)->next;
8.    while((*r)->key < k){
9.
10.     l = *r;
11.     *r = r_next;
12.
13.     r_next = (*r)->next;
14.   }
15.}
```

# Concurrent set – Attempt 1

- PESSIMISTIC approach
- Synchronize via global lock

INSERT(35)

INSERT(25)

DELETE(40)

INSERT(55)

| H | | 10 | | 20 | | 30 | | 40 | | 50 | | | |

# Concurrent set – Attempt 1 (SRC)

```
1. bool do_operation(int k, int op_type){
2.    bool res = true;
3.    node *l,*r;
4.    LOCK(&glock);
5.    l = search(k, &r);
6.    switch(op_type){
7.      case(INSERT):
8.        if(r->key == k)
9.          res = false;
10.       else
11.         l->next = new node(k,r);
12.       break;
13.     case(DELETE):
14.       if(r->key == k)
15.         l->next = r->next;
16.       else
17.         res = false;
18.       break;
19.   }
20.   UNLOCK(&glock);
21.
22.   return res;
23.}
```

```
1. node* search(int k, node **r){
2.    node *l, *r_next;
3.    l = set->head;
4.
5.    *r = l->next;
6.
7.    r_next = (*r)->next;
8.    while((*r)->key < k){
9.
10.     l = *r;
11.     *r = r_next;
12.
13.     r_next = (*r)->next;
14.   }
15.}
```
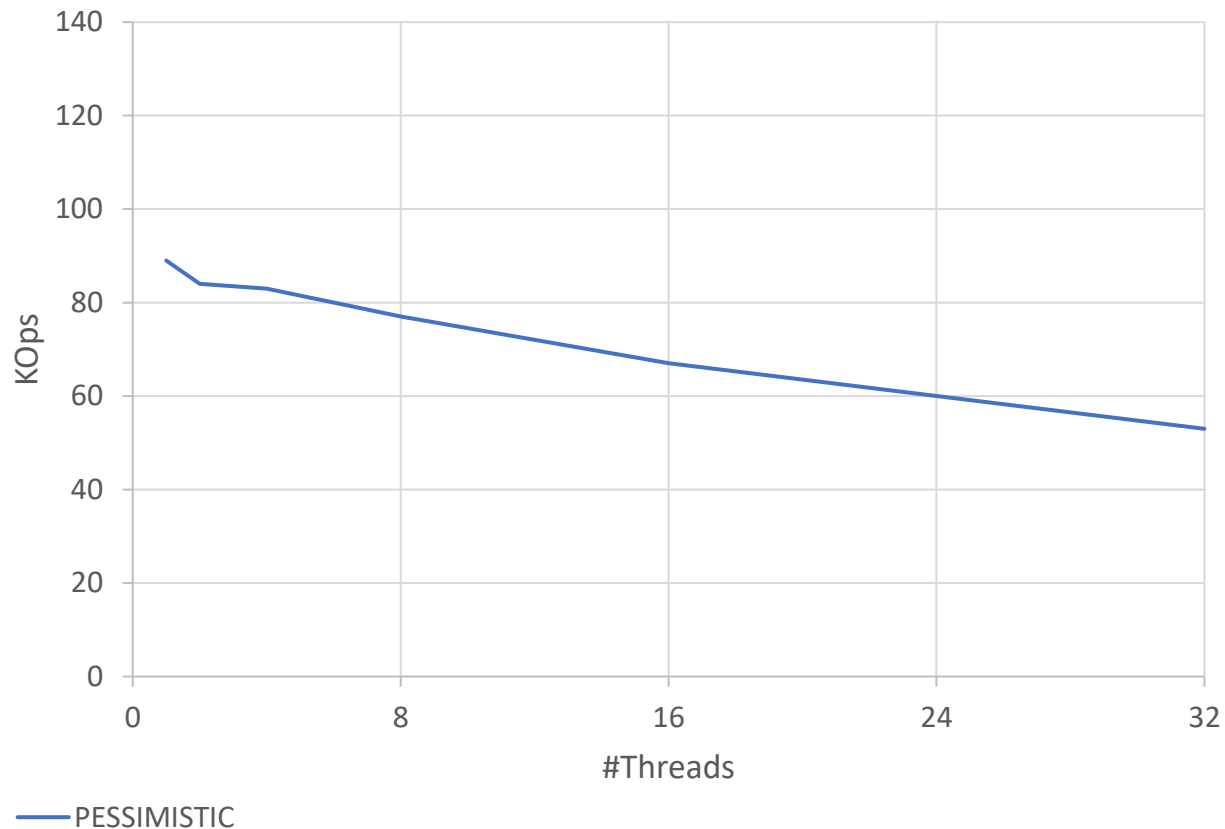
# Concurrent set – Attempt 1

AMD Opteron 6128 – 32Cores

KeyRange = [0,6000]          SetSize = 2400          Update=100%

# Concurrent set – Attempt 1



INSERT(5)

H — 10 — 20 — 30 — 40 — 50 —

# Concurrent set – Attempt 1

- PESSIMISTIC approach
- Synchronize via global lock
⇒ NO SCALABILITY!

# Concurrent set – Attempt 2

- Fine-grain approach
- Each node has its own lock
- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

# Search algorithm

INSERT(55)

H — 10 — 20 — 30 — 40 — 50 — T

# Search algorithm

- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

# Search algorithm

- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

# Search algorithm

- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

# Search algorithm

- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

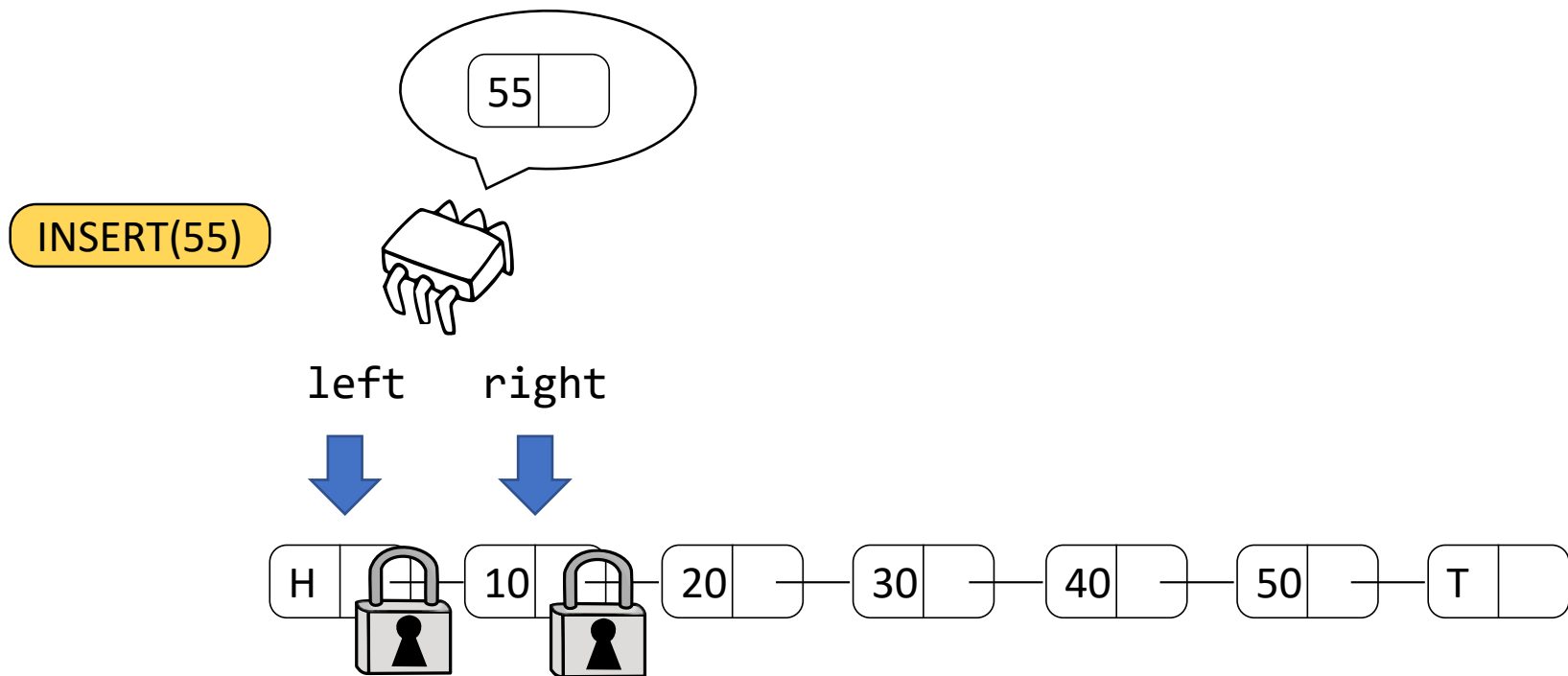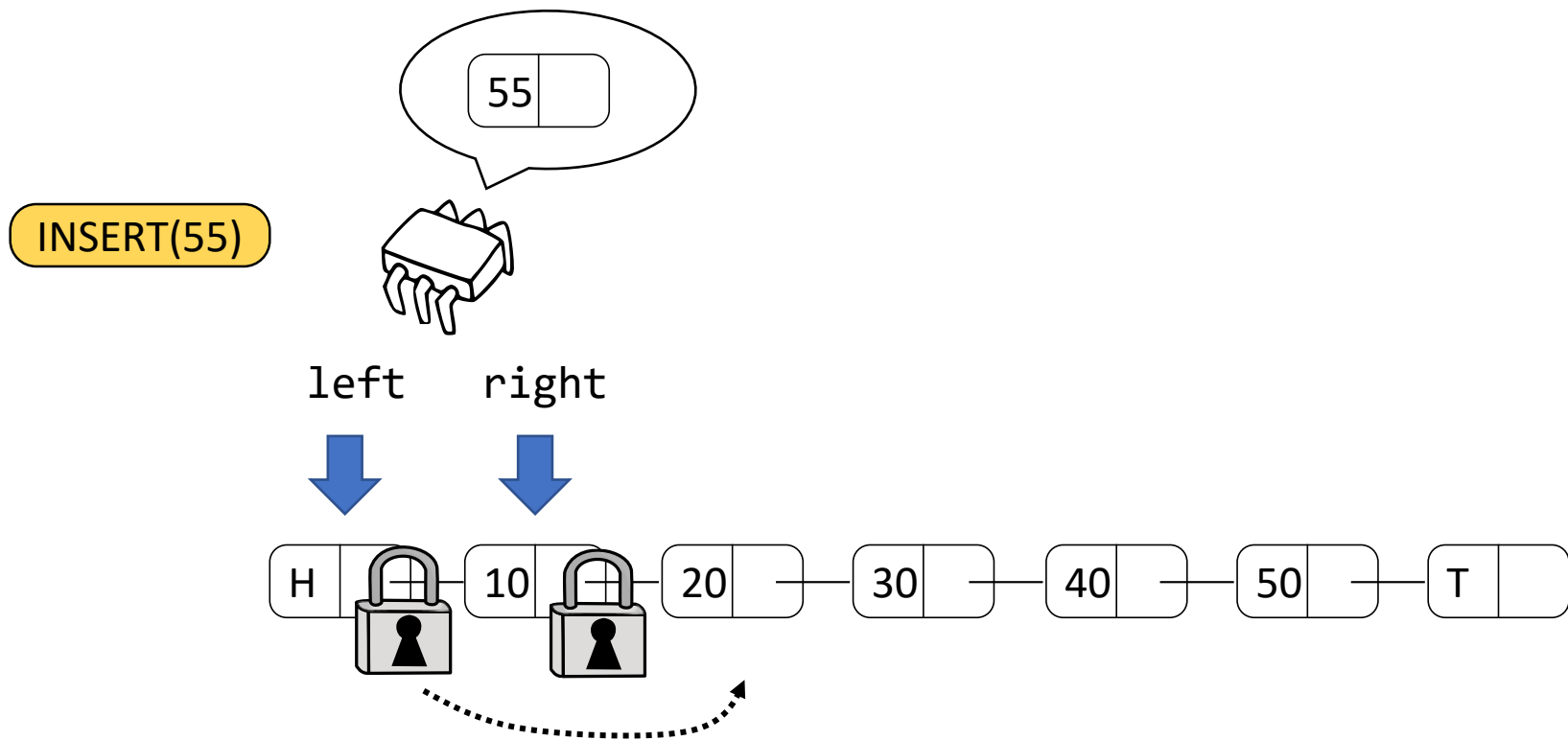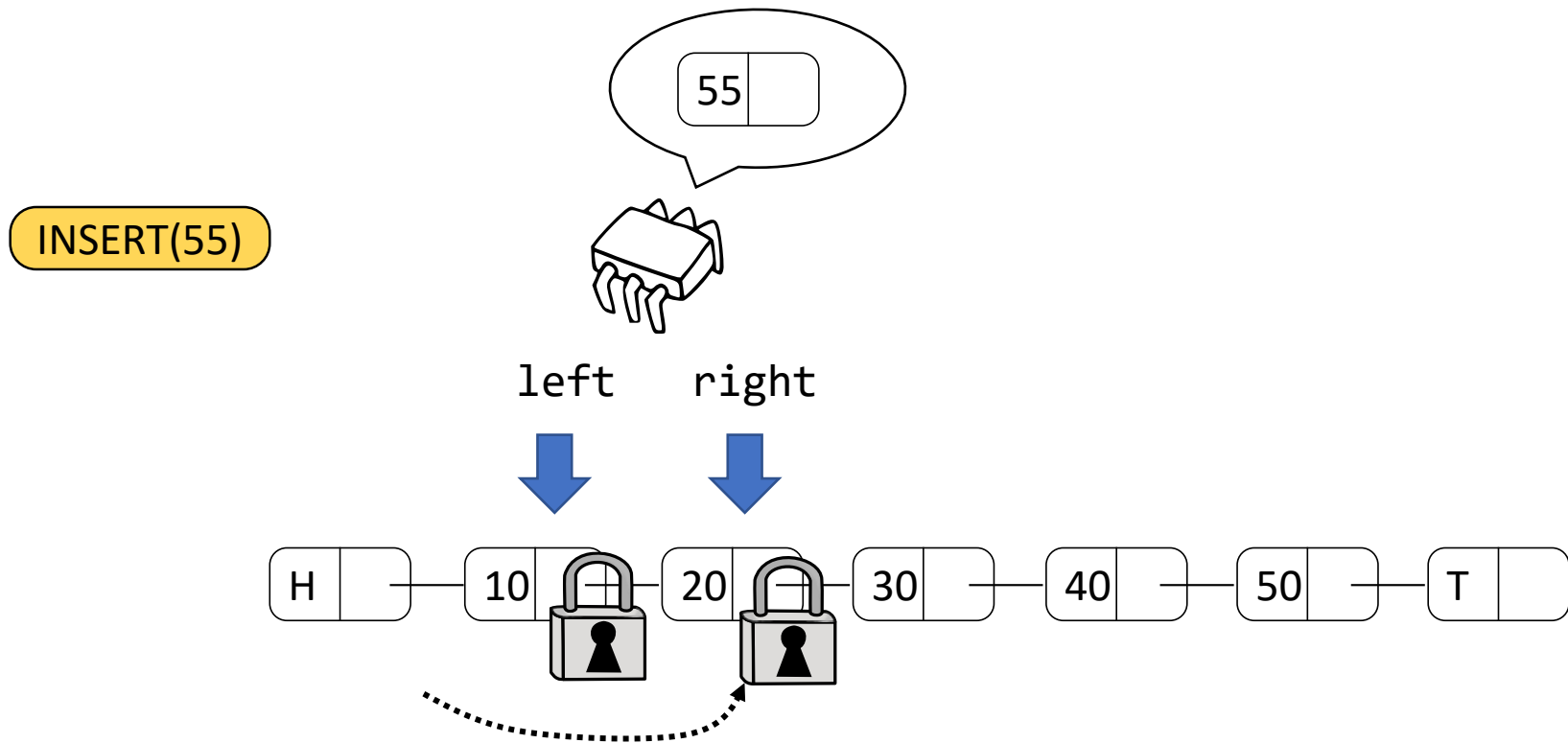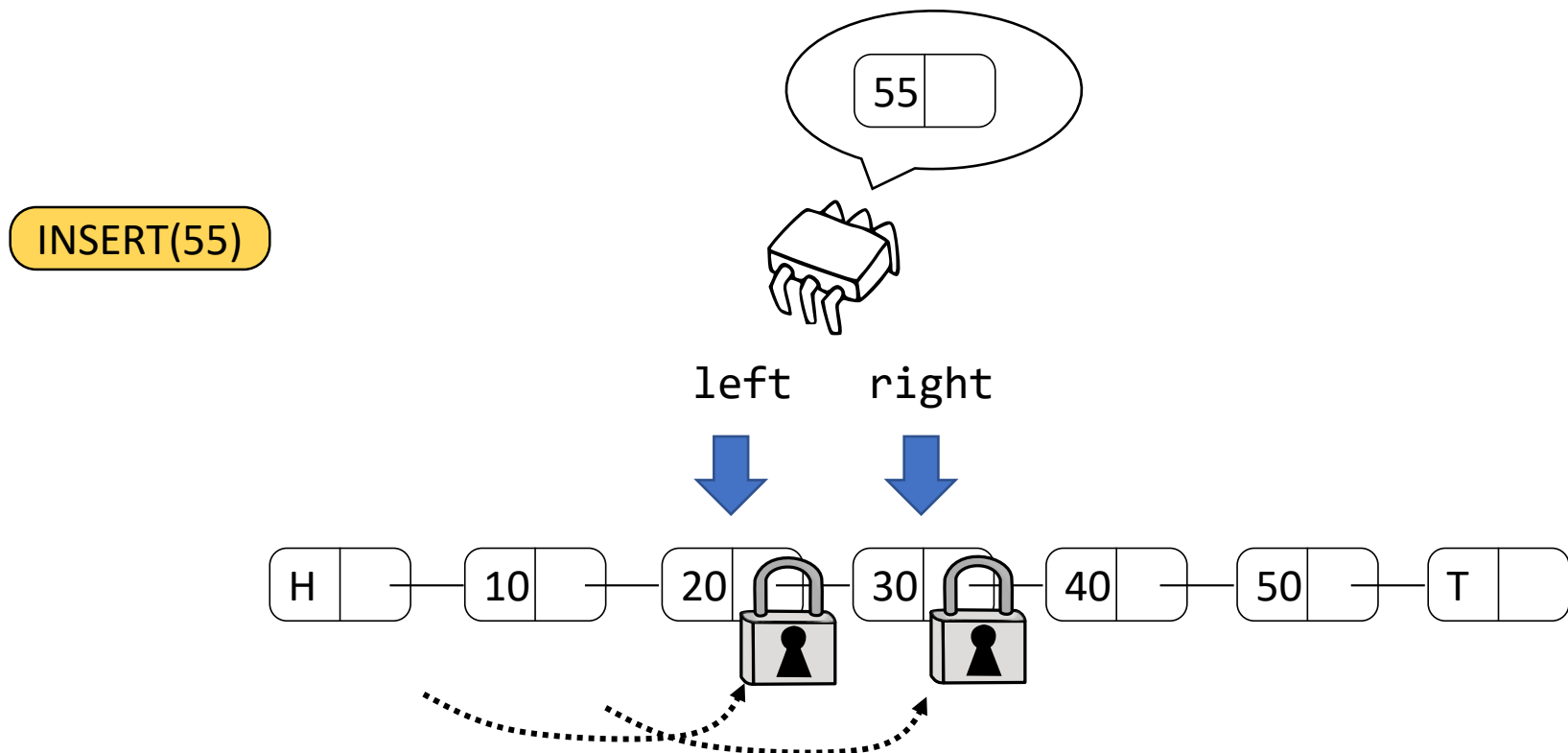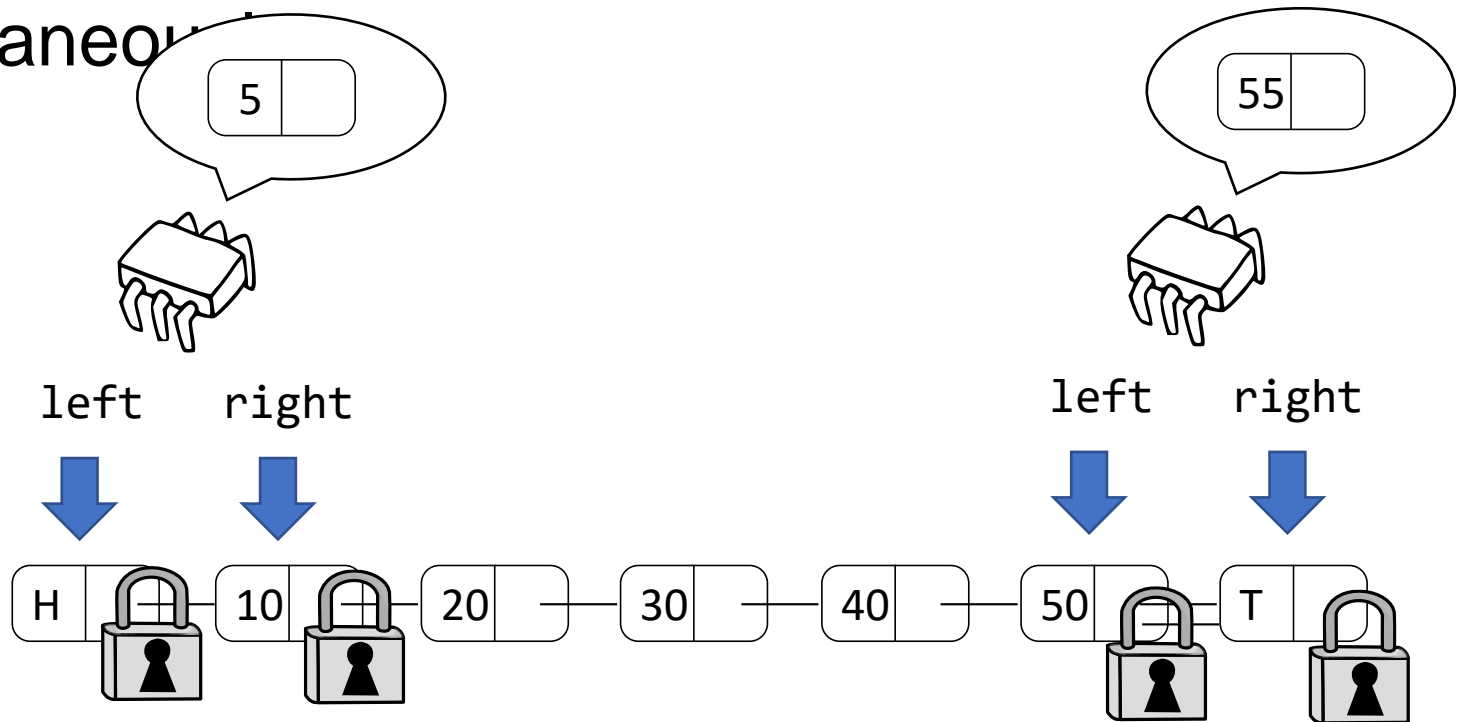# Search algorithm

- Keep two locks at a time (lock coupling):
  - One on the current node
  - One on its predecessor

- Multiple threads access the data structure simultaneously

# Concurrent set – Attempt 2 (SRC)

```
1.  bool do_operation(int k, int op_type){
2.     bool res = true;
3.     node *l,*r;
4.     LOCK(&glock);
5.     l = search(k, &r);
6.     switch(op_type){
7.       case(INSERT):
8.         if(r->key == k)
9.            res = false;
10.        else
11.           l->next = new node(k,r);
12.        break;
13.      case(DELETE):
14.        if(r->key == k)
15.           l->next = r->next;
16.        else
17.           res = false;
18.        break;
19.     }
20.     UNLOCK(&glock);
21.     UNLOCK(&l->lock);
22.     UNLOCK(&r->lock);
23.     return res;
24. }
```

```
1.  node* search(int k, node **r){
2.     node *l, *r_next;
3.     l = set->head;
4.     LOCK(&l->lock);
5.     *r = l->next;
6.     LOCK(&(*r)->lock);
7.     r_next = (*r)->next;
8.     while((*r)->key < k){
9.        UNLOCK(&l->lock);
10.       l = *r;
11.       *r = r_next;
12.       LOCK(&(*r)->lock);
13.       r_next = (*r)->next;
14.    }
15. }
```

# Concurrent set – Attempt 2

AMD Opteron 6128 – 32Cores

KeyRange = [0,6000]          SetSize = 2400          Update=100%

# Search algorithm

• Allows an increased parallelism but…

# Search algorithm

- Allows an increased parallelism but…
- High costs for lock handover

# Recap

- Explored two <u>blocking</u> strategies:
1. Global (coarse-grain) lock



2. (Fine-grain) Lock coupling

# Concurrent set – Attempt 3

DELETE(40)

INSERT(55)

| H | | 10 | | 20 | | 30 | | 40 | | 50 | | T | |

# Concurrent set – Attempt 3

- NON-BLOCKING approach [Harris linked list]
- Search without acquiring any lock
- Apply updates with individual atomic instructions

DELETE(40)

INSERT(55)

H → 10 → 20 → 30 → 40 → 50 → T

# Non-blocking insert & delete algorithms

Insert:

1. Search left and right nodes

2. Insert the new item with a CAS

3. If CAS fails restart from 1

Delete:

1. Search left and right nodes

2. Disconnect the item with a CAS

3. If CAS fails restart from 1



INSERT(20)

DELETE(10)

- Is it correct?

# Incorrect delete algorithm

- Edge cases might lead to losing items!

# Incorrect delete algorithm

- Edge cases might lead to losing items!



1. Thread A gets left and right node and go to sleep
2. Thread B disconnects the node containing 10
3. Thread A wakes up and add 20 after 10
4. The new item is lost

# Incorrect delete algorithm

- Edge cases might lead to losing items!



1. Thread A gets left and right node and go to sleep
2. Thread B disconnects the node containing 10
3. Thread A wakes up and add 20 after 10
4. The new item is lost

# Incorrect delete algorithm

- Edge cases might lead to losing items!



1. Thread A gets left and right node and go to sleep
2. Thread B disconnects the node containing 10
3. Thread A wakes up and add 20 after 10
4. The new item is lost

# Incorrect delete algorithm

- Edge cases might lead to losing items!



1. Thread A gets left and right node and go to sleep
2. Thread B disconnects the node containing 10
3. Thread A wakes up and add 20 after 10
4. The new item is lost

# The correct delete algorithm

- Adopt logical deletion:

1. Get left and right node

2. Mark the item as deleted via CAS (*logical* deletion)

3. If CAS fails GOTO 1

4. Disconnect the item via CAS (*physical* deletion)

5. If CAS fails GOTO 4

# The correct delete algorithm

- Adopt logical deletion

1. Ge...

2. Ma... *logical* de...

3. If ...

4. Di... *cal* deletion)

5. If CAS fails GOTO 4

**CAS**

| key | next |
|-----|------|
| 10 | 0xff ... ✗1 mark |

- Typically memory objects are byte aligned
- The LSB is always 0! BIT STEALING!!!

# The correct delete algorithm

DELETE(10)



INSERT(20)

left          right

- Updates of the "next" field by two opposite concurrent operations cannot both succeed
- What to do upon conflict (failed CAS)?   RESTART FROM SCRATCH!!

# Non-blocking search

- The search returns two adjacent <u>non-marked</u> (left and right) nodes

- Housekeeping: disconnect logically delete items during searches

# Non-blocking search

- The search returns two adjacent <u>non-marked</u> (left and right) nodes

- Housekeeping: disconnect logically delete items during searches

# Non-blocking search

- The search returns two adjacent <u>non-marked</u> (left and right) nodes

- Housekeeping: disconnect logically delete items during searches

# Non-blocking search

- The search returns two adjacent <u>non-marked</u> (left and right) nodes

- Housekeeping: disconnect logically delete items during searches

# Concurrent set – Attempt 3 (SRC)

```
1. bool do_operation(int k, int op_type){
2.    node *l,*r, *n = new node(k);
3.    l = search(k, &r);                /* get left and right node */
4.    switch(op_type){
5.      case(INSERT):
6.        if(r->key == k) return false;  /* key present in the set */
7.        n->next = r;
8.        l->next = n;                   /* insert the item      */
9.
10.
11.       break;
12.     case(DELETE):
13.       if(r->key != k) return false;  /* key not present      */
14.       l->next = r->next;             /* remove the key       */
15.
16.
17.
18.       break;
19.    }
20.    return true;
21.}
```

# Concurrent set – Attempt 3 (SRC)

```
1. bool do_operation(int k, int op_type){
2.    node *l,*r, *n = new node(k);
3.    l = search(k, &r);                    /* get left and right node */
4.    switch(op_type){
5.      case(INSERT):
6.        if(r->key == k) return false;  /* key present in the set */
7.        n->next = r;
8.        l->next = n;                       /* insert the item      */
9.        if(!CAS(&l->next, r, n))
10.          goto 3;         /* insertion failed the item -> restart */
11.       break;
12.     case(DELETE):
13.        if(r->key != k) return false;  /* key not present       */
14.        l->next = r->next;                /* remove the key        */
15.        if(is_marked_ref((l=r->next)) || !CAS(&r->next, l, mark(l)))
16.          goto 3;         /* insertion failed the item -> restart */
17.       search(k,&r);                       /* repeat search        */
18.       break;
19.   }
20.   return true;
21.}
```

# Concurrent set – Attempt 3 (SRC)

```
1.  node* search(int k, node **r){
2.     node *l, *t, *t_next, *l_next;
3.     *t = set->head;
4.      t_next = t->head->next;
5.      while(1){                        /* FIND LEFT AND RIGHT NODE */
6.          if(!is_marked(t_next)){
7.              l = t;
8.              l_next = t_next;
9.          }
10.         t = get_unmarked_ref((t_next);
11.         t_next = t->next;
12.         if(!is_marked_ref(t_next) && t->key >= k) break;
13.     }
14.     *r = t;
15.     /* DEL MARKED NODES */
16.     if(l_next != *r && !CAS(&l->next, l_next, *r) goto 3;
17.     return l;
18. }
```

# Concurrent set – Attempt 3

AMD Opteron 6128 – 32Cores

KeyRange = [0,6000]          SetSize = 2400          Update=100%



Legend: PESSIMISTIC, CHAINED, LOCK-FREE

# Safety and liveness guarantees

- The algorithm is NON-BLOCKING (LOCK-FREE):
    - If a thread A is stuck in a retry loop => a CAS fails each time
    - If a CAS fail, it is because of another CAS that was successfully executed by a thread B
    - Thread B is making progress

- The algorithm is LINEARIZABLE:
    - Each method execution take effect in an atomic point (a successful CAS) contained between its invocation and reply
    - The order obtained by using these points has been proved to be compliant with the Set semantic

# Progress (Lock freedom)

- Each method update method has two main steps
  - A search, which might end with a CAS
  - A CAS to insert delete a node

1. Suppose an update method is stuck in a search:
   - The key range is finite, so the number of node is finite
   - It continuously fails to disconnect marked nodes
   - It means that new nodes have been both inserted and marked!
     - Other threads have completed update methods

2. Suppose an updated method always fails its last step (insertion or marking)
   - Other threads have modified the target next pointer
   - If it is due to the disconnection of marked nodes, see point 1
   - If it is due to the updated step other methods have completed

# Safety (Linearizability)

1. The search returns 2 adjacent nodes in an atomic point
   1. The read of next field of the left node
   2. The CAS that make left and right adjacent

- It is like that the search made a snapshot of interested key interval

2. Find, unsuccessful delete and unsuccessful insert linearize with the search (1.1 or 1.2)

3. Insert linearizes with the successful CAS to connect a new node

4. Delete linearizes with the successful CAS to mark a node

# Problems

- It is not possible to flip a bit of a reference on memory-managed languages (e.g. JAVA)

- How to solve?

# Locks + Optimism

- Use one lock per node
- Move "marked" to a dedicated field

# Locks + Optimism (insert)

- Use one lock per node
- Move "marked" to a dedicated field
- Find left and right without taking locks!
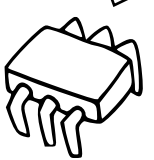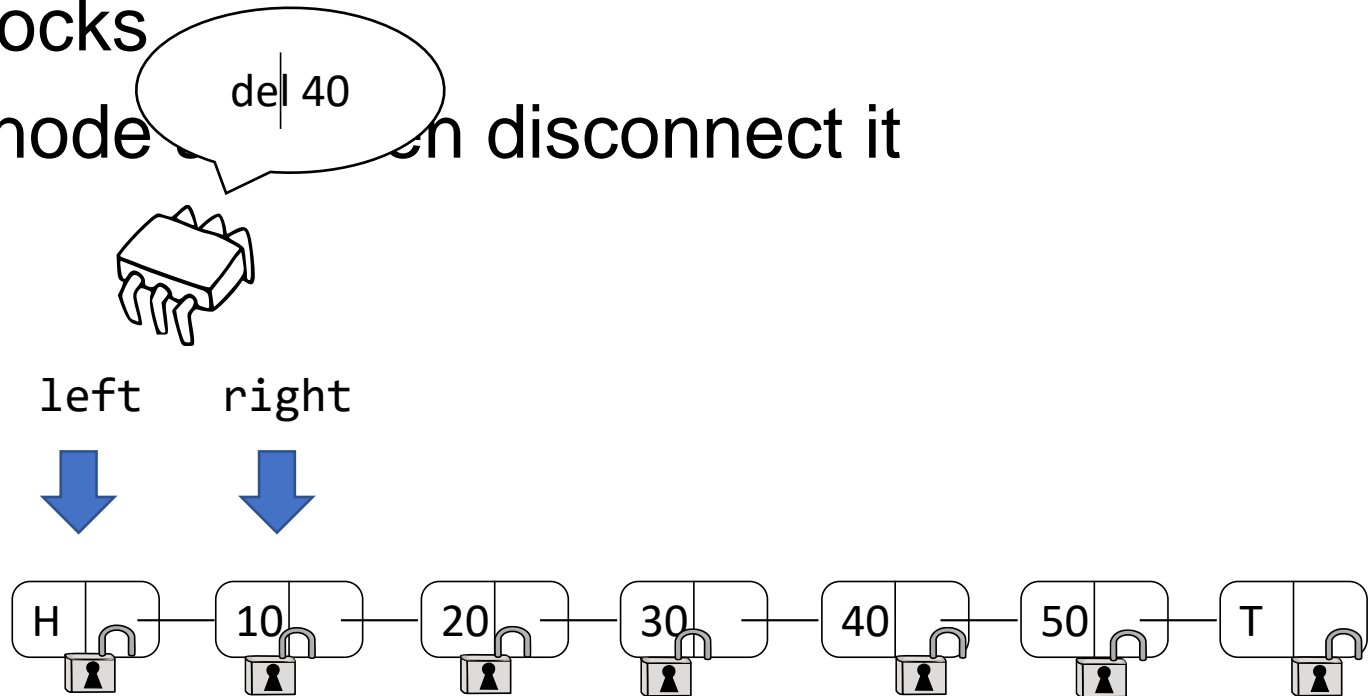- Take locks
- Insert the n

# Locks + Optimism (delete)

- Use one lock per node
- Move "marked" to a dedicated field
- Find left and right without taking locks!
- Take locks
- Mark node and then disconnect it

del 40

left    right

H 10 20 30 40 50 T

# Locks + Optimism (delete)

- Why "optimistic"? Do work (search) and hope nothing wrong happens!

- What could go wrong?



left    right

# Locks + Optimism (delete)

- Why "optimistic"?  Do work (search) and hope nothing wrong happens!

- What could go wrong?
  - Left and/or right being marked
  - Left and right not adjacent

- How to solve?

- Validation of search results:
  - Left unmarked
  - Right unmarked
  - Left.next = right

# Locks + Optimism (delete)

- Why "optimistic"?  Do work (search) and hope nothing wrong happens!

- What could go wrong?
    - Left and/or right being marked
    - Left and right not adjacent

- How to solve?

- Validation of search results:
    - Left unmarked
    - Right unmarked
    - Left.next = right

# Locks + Optimism = Lazy List

- What about correctness?
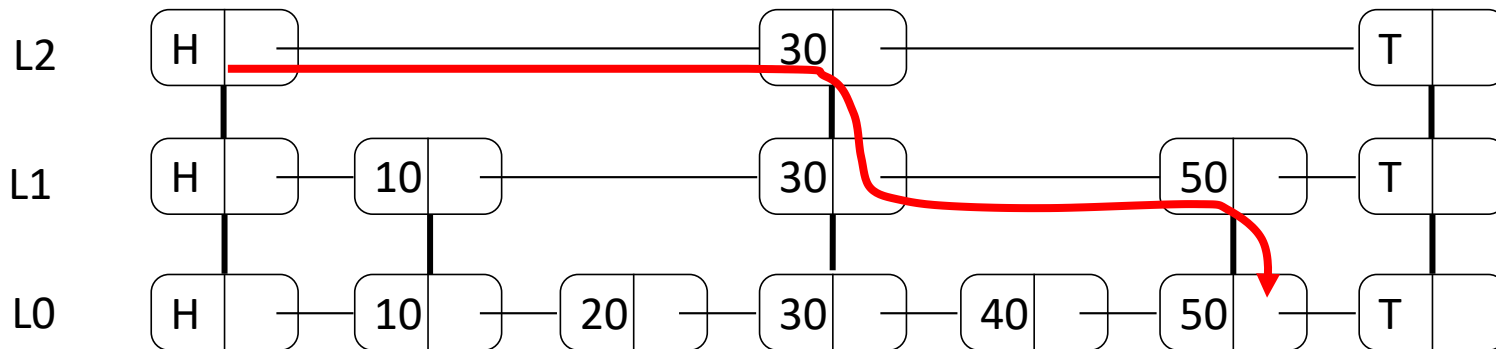- What about progress?

# Can we do better?

- Costs: O(n)
- Starting from scalable "simple" set implementation we can build faster set implementations
  - Hash table: O(1)
    - Array of buckets
    - Buckets are concurrent ordered-list based sets
- We know that a search in an ordered set could be more efficient $O(\log(n))$
- How?

# Skip list [Pugh 1990]

- Generalization of sorted linked lists

- Randomized data structure

- Costs: O(log(n))


- Idea:
  1. Maintain a core sorted linked list L0
  2. Use additional sorted linked lists Li such that:
     1. Li $\subset$ Li-1
     2. |Li| $\approx$ |Li-1|/2
  3. Searches use lists in decreasing order
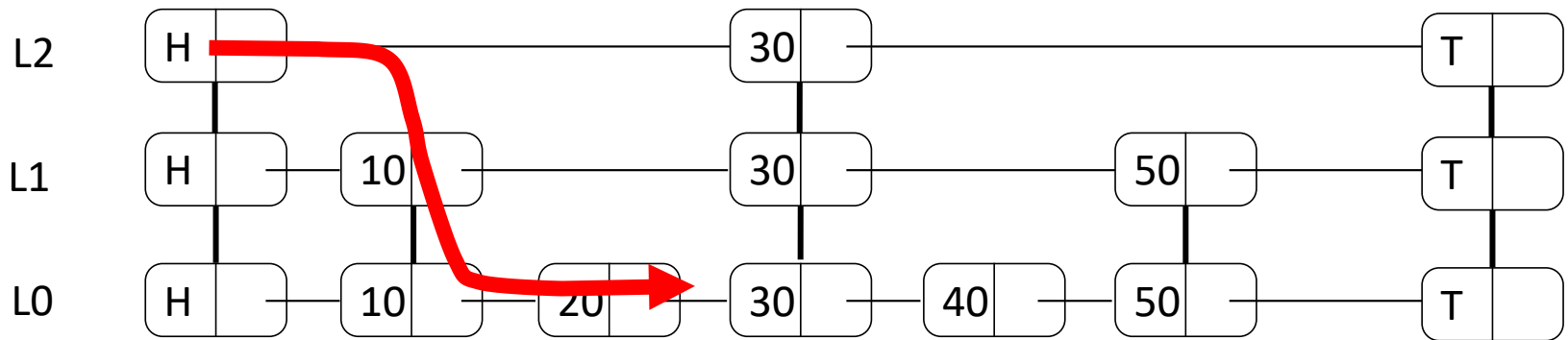
# Skip list [Pugh 1990]

Search(50)

L2  H — 30 — T

L1  H — 10 — 30 — 50 — T

L0  H — 10 — 20 — 30 — 40 — 50 — T

# Skip list [Pugh 1990]

Insert(25)

# Skip list [Pugh 1990]

Insert(25)

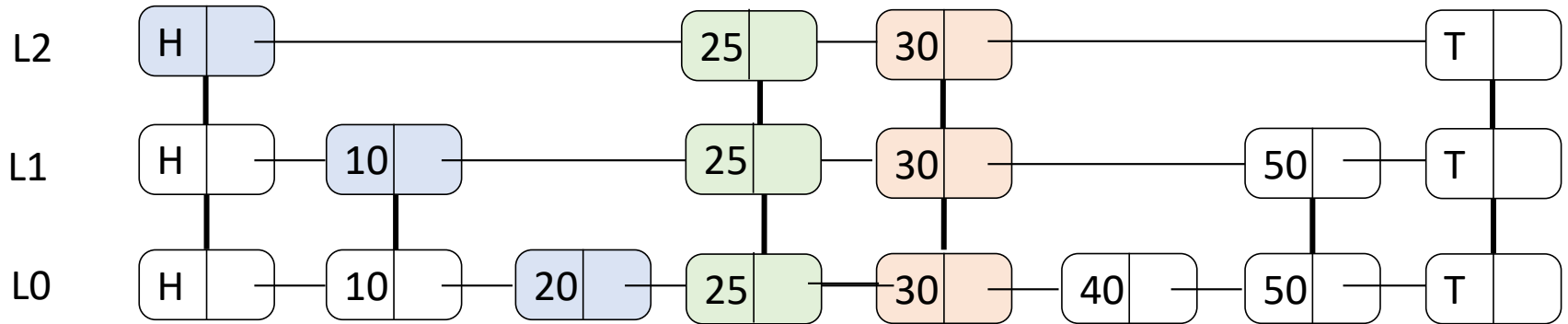| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| L2 | H | | | 25 | 30 | | | T |
| L1 | H | 10 | | 25 | 30 | | 50 | T |
| L0 | H | 10 | 20 | 25 | 30 | 40 | 50 | T |

Should I insert 25 at L1?  Flip a coin!

Should I insert 25 at L2?  Flip a coin!

# Skip list [Pugh 1990]

Delete(25)

# Skip list [Pugh 1990]
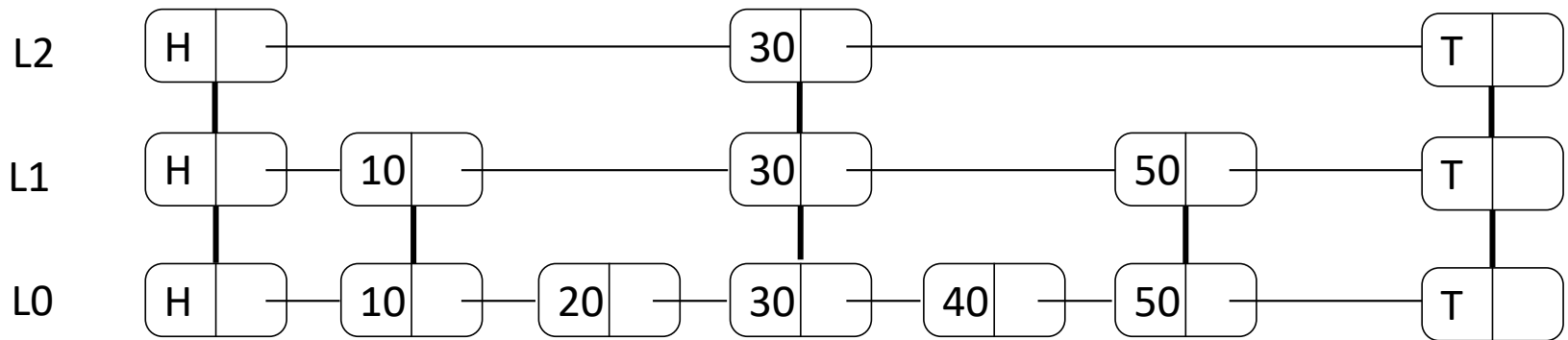
- How many (expected) keys for each level?
- L0 = N
- L1 = N/2
- L2 = N/4
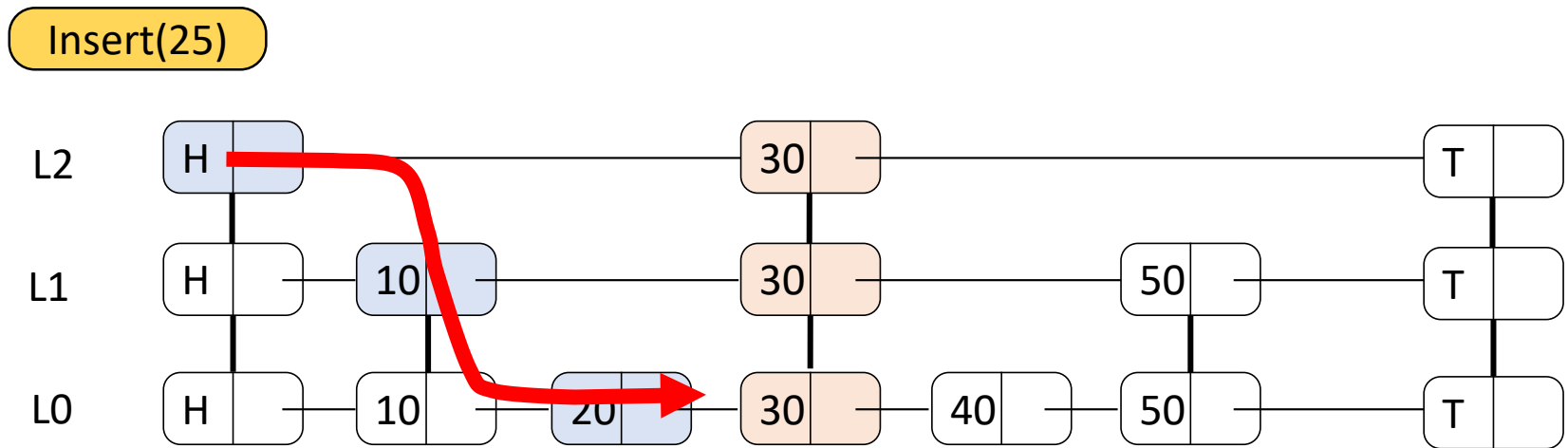
…

- L(logN) = 1

# Skip list [Pugh 1990]

- How many steps per level?

# Non-blocking Skip list [Fraser2004]

Insert(25)
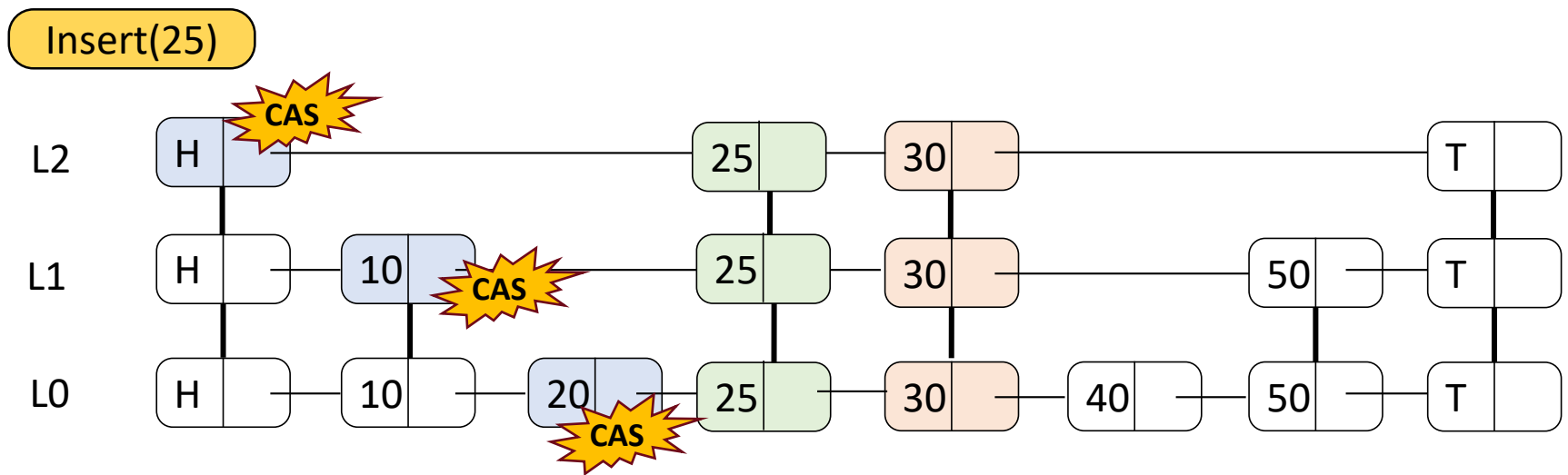
# Non-blocking Skip list [Fraser2004]

Insert(25)

# Non-blocking Skip list [Fraser2004]

Insert(25)

# Non-blocking Skip list [Fraser2004]

Delete(25)

# Non-blocking Skip list [Fraser2004]

# Concurrent Data Structures:
# Priority queues

# Priority queue implementations

- Priority Queue methods:
  - `enqueue(k):` adds a new item
  - `dequeue():` returns and remove the highest priority item

- Implemented as an ordered linked list ←

This is a huge simplification.
Tipically they are implemented as
skip-lists (log(n)) or calendar queues
(O(1))

ENQ(35)

ENQ(25)

DEQ()

ENQ(55)

H → 0.1 → 1.3 → 5.0 → 6.5 → 7.1 → 9.8 → T

# Priority queue – Attempt 1

- Enqueue: works as insertions in the non-blocking Set
  - Connect via CAS

- Dequeues: work as deletions in the non-blocking Set
  - Mark as logically deleted, but
  - DISCONNECT IMMEDIATELY

- Is it scalable?

# Priority queue – Attempt 1



Queue Size ≈ 256000

Scalability collapse

SLCQ △
NBCQ ○

Throughput (MOps/s) vs #Threads

# Priority queues: an inherently "sequential" semantic

- Enqueue offers a high level of disjoint access parallelism
- Dequeues are prone to conflicts



This region is highly shared among processors' caches

# Lazy deletion within priority queues

- If we use lazy deletion "as is", we might obtain non-linearizable extractions



Non-linearizable extraction

# Correct lazy deletion within priority queues

- To implement correct extractions with lazy deletions there are two main approaches

1. Move the logical mark of a node in the field "next" of its predecessor

# Correct lazy deletion within priority queues

- To implement correct extractions with lazy deletions there are two main approaches

2. Use logical timestamps:
   - incremented each time a new minimum has been inserted
   - extract item compatible with the timestamp read at the beginning

Ts=0

!

| H | → | 0.1 | → | ✖ | → | ✖ | → | 6.5 | → | 7.1 | → | 9.8 | → | T |

~~Ts=0~~       Ts=0        Ts=0        Ts=0        Ts=1        Ts=0        Ts=0

Ts=1

# PQ – Attempt 2 - Introducing Conflict Resiliency

- Lazy deletion
- Skip logically deleted items $\Rightarrow$ IT INCREASES THE NUMBER OF STEPS $\Rightarrow$ EXPENSIVE IN TERMS OF IMPACT ON CACHE
- Periodic Housekeeping

#deleted items> threshold

# Priority queue – Attempt 2



Queue Size ≈ 256000

# On the conflict resiliency trade off

- The number of steps per dequeue and costs of housekeeping are **<u>dependent</u>**:

⬆ **THRESHOLD** ⟹ ⬆ **READ LATENCY** ☹ **and** ⬇ **RMW IMPACT** ☺

⬇ **THRESHOLD** ⟹ ⬇ **READ LATENCY** ☺ **and** ⬆ **RMW IMPACT** ☹

# Conflict resiliency trade offs



Queue Size ≈ 2560000

# Priority queues – Attempt 3



Queue Size ≈ 2560000

# Open challenges

How to achieve scalability for priority queues?

- NO ANSWER for correct priority queue

- The research moved on looking for RELAXED SEMANTICS for priority queues
  - Enable scalability for extractions by removing an item which is "near" the minimum

- Explore orthogonal approaches by guaranteeing RELAXED CORRECTNESS CONDITIONS
  - K-linearizability
  - Quasi-linearizabilty
  - Quiescent consistency
  - Sequential consistency?

- Explore new hardware capabilities (e.g. HTM)

# Why linearizable non-blocking algorithms?

- Performance is a good reason, but not the unique one

- The composition of linearizable algorithm is still linearizable

- Blocking algorithms (and their composition) might suffer from deadlocks, priority inversions and convoying

- The composition of non-blocking algorithms is non-blocking as a whole (progress property of individual algorithm might be hampered)

- Libraries should implement their algorithms in a non-blocking linearizable fashion
  - E.g. Java implements non-blocking concurrent data structure

# Concurrent Data Structures:
# FIFO queues

# FIFO queue implementation

- Queue methods:
  - `enqueue(v)`
  - `dequeue()`

- Implemented as a linked list

# FIFO queue implementation

- Slightly different
- One dummy node, two pointers to access the data structure:
  - Head: points ALWAYS to a DUMMY node item
  - Tail: SHOULD point to the youngest item

# FIFO queue implementation

- Insert:

1. Get node pointed by tail

2. Scan until next is NULL

3. Try to insert with CAS

4. If KO goto 1

5. Else try to update Tail

- Dequeue:

1. Get node pointed by head

2. Try to update head with its next

3. If KO goto 1

ENQ(c)

DEQ()

H

T

CAS

c

NULL

DU

a

b

CAS

NULL

H

CAS

T

DU

a

b

NULL

This becomes the new dummy node

# The whole story

- Since the insertion of a new item and the tail update are two separate RMW they might be inconsistent

- Also dequeuers might need to update tail before updating head

- This ensures that TAIL cannot go behind HEAD

# Emptiness condition

- There is a NULL node after the one pointed by HEAD

DEQ()

H

T

DU → a → b → NULL

# Wait-free FIFO queue

- What about a wait-free queue?

- Wait-free means that all method invocations are guaranteed to complete

- Can we modify the lock-free FIFO queue to achieve this?

- Lock-free means that some thread might starve

- If before starting any new operation we complete a pending operation, all method invocation complete eventually

# Wait-free FIFO queue

- We need to be aware of pending calls

| phase |
|-------|
| Pending |
| isEnqueue |
| Node |

| 9 |
|-------|
| True |
| False |
| NULL |

| 4 |
|-------|
| False |
| True |
| NULL |

| 9 |
|-------|
| False |
| True |
| NULL |

- Split operations on the linked list into 2 steps:
  1. Modify nodes for enqueue/dequeue (main step)
  2. Modify head/tail pointers (finishing step)

# Wait-free FIFO queue

- Enqueue/Dequeue structure
    1. Publish op record
    2. Get the set **S** of all pending ops whose record has been previously or concurrently published
    3. Help any operation in **S**
    4. Do a finishing step

# Wait-free FIFO queue

- Enqueue/Dequeue structure
    1. Publish op record

# Wait-free FIFO queue

- Enqueue/Dequeue structure
    2. Get the set **S** of all pending ops whose record has been previously or concurrently published

# Wait-free FIFO queue

- Enqueue/Dequeue structure
    3. Help any operation in **S** (dequeue)
        a. Main step

# Wait-free FIFO queue

- Enqueue/Dequeue structure
  3. Help any operation in **S** (dequeue)
     a. Main step
     b. Finishing step

# Wait-free FIFO queue

- Enqueue/Dequeue structure
  3. Help any operation in **S** (enqueue)
     a. Main step

# Wait-free FIFO queue

- Enqueue/Dequeue structure
  3. Help any operation in **S** (enqueue)
     a. Main step
     b. Finishing step

# Wait-free FIFO queue

- Enqueue/Dequeue structure
    1. Publish op record
    2. Get the set **S** of all pending ops whose record has been previously or concurrently published
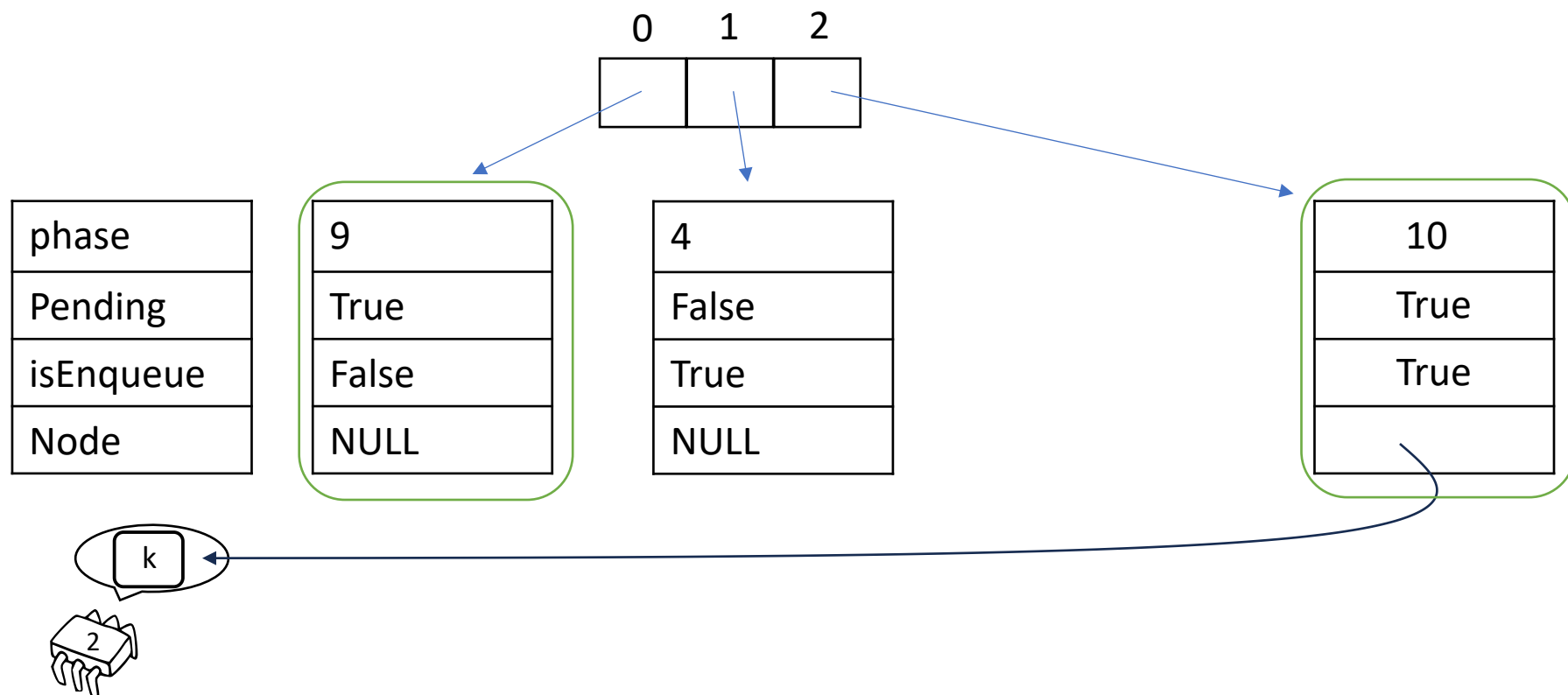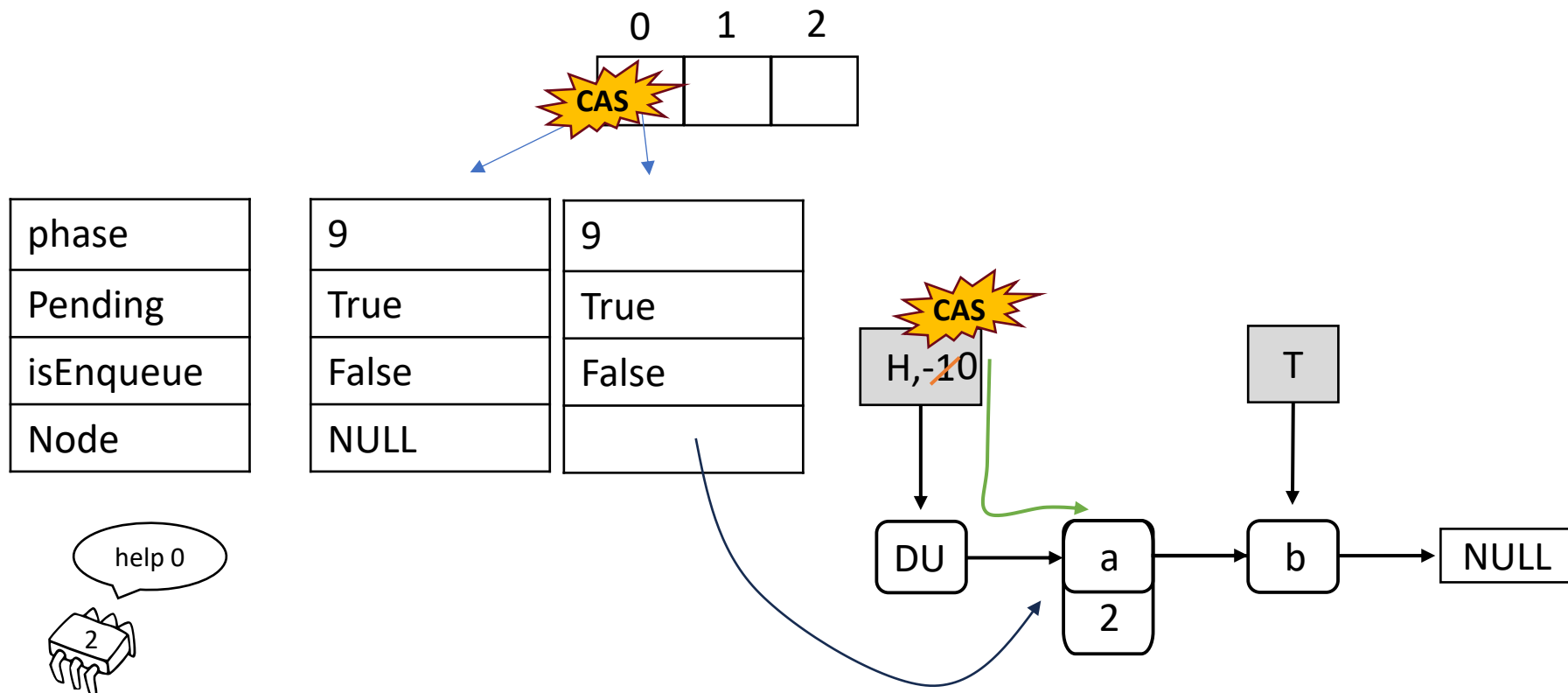    3. Help any operation in **S**
    4. Do a finishing step

Opt 1: help only one pending op
Opt 2: use FAD to get phase num.



total completion time

# Fast Wait-free FIFO queue

- Try with lock-free approach
- If starving, back-off to wait-free implementation

# Concurrent Data Structures:

# Atomic MRSW registers

# Atomic MRSW Register

# Atomic MRSW Register



WRITER → VALUE → READER $_1$ ... READER $_i$ ... READER $_N$

# Atomic MRSW Register - Peterson

| Variable | Type | Description |
|---|---|---|
| BUF1 | Buffer | Main buffer. |
| BUF2 | Buffer | Backup buffer. |
| COPYBUF | Array of $n$ buffers | An individual backup buffer for each reader. |

**Write operation:**

```
PW1    WFLAG := true;
PW2    write to BUF1;
PW3    SWITCH := !SWITCH;
PW4    WFLAG := false;
PW5    for (each reader r)
PW6        if (READING[r] != WRITING[r])
PW7            write to COPYBUF[r];
PW8            WRITING[r] := READING[r];
PW9    write to BUF2;
```

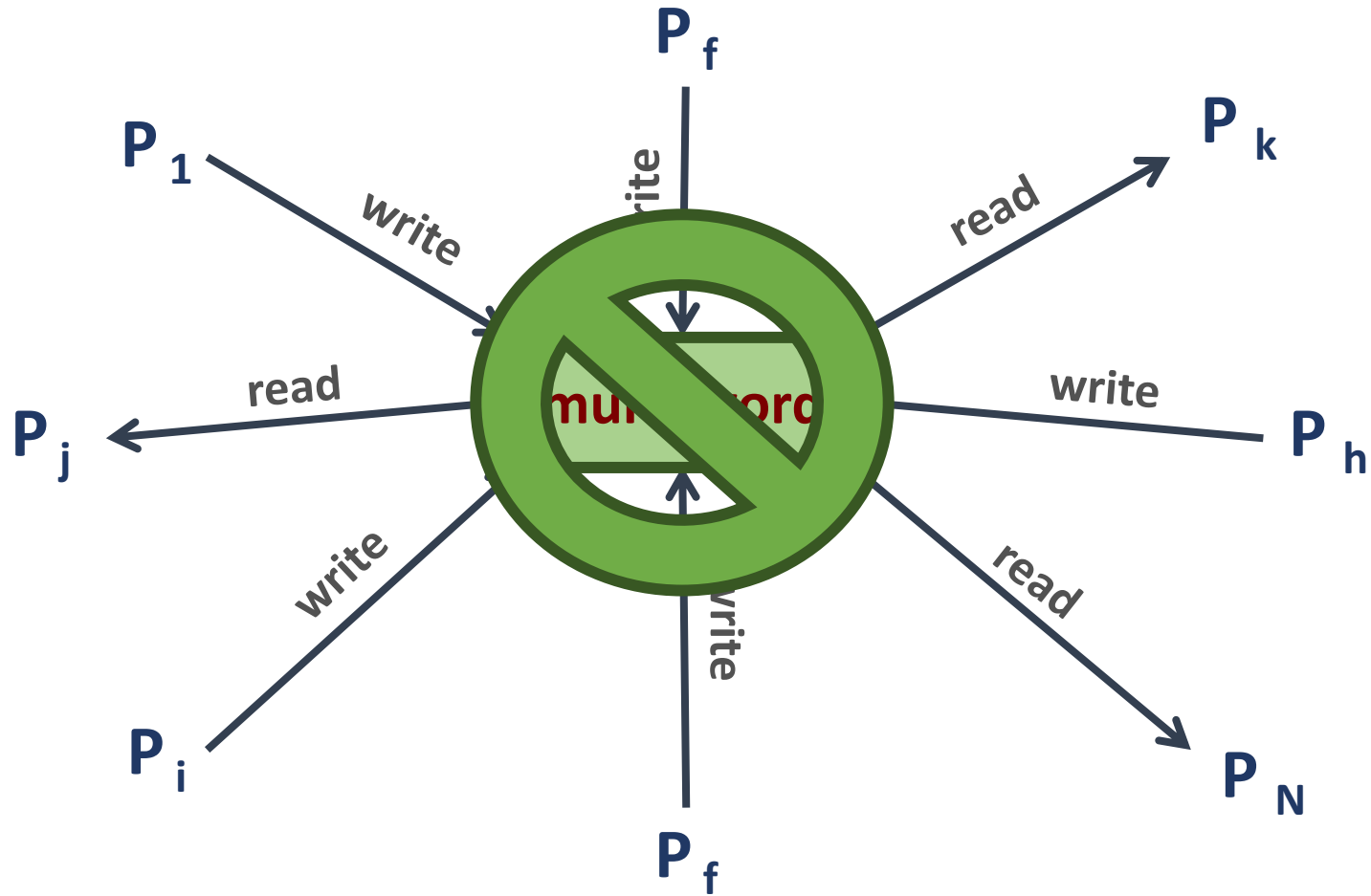**Read operation by reader r:**

```
PR1    READING[r] := !WRITING[r];
PR2    flag1 := WFLAG;
PR3    sw1 := SWITCH;
PR4    read BUF1;
PR5    flag2 := WFLAG;
PR6    sw2 := SWITCH;
PR7    read BUF2;
PR8    if (READING[r] == WRITING[r])
PR9        return the value in COPYBUF[r];
PR10   else if ((sw1 != sw2) || flag1 || flag2)
PR11       return the value read from BUF2;
PR12   else
PR13       return the value read from BUF1;
```

# Atomic MRSW Register - Larsson

| Variable | Type | Description |
|---|---|---|
| BUF$[n+2]$ | Array of $n+2$ buffers | The buffers for the register value. |

**Read operation by reader r:**

R1      readerbit := 1 << (r + PTRFIELDLEN);

R2      rsync := fetch_and_or(**&SYNC**, readerbit);

R3      rptr := rsync **&** PTRFIELD;

R4      **read** BUF[rptr]

**Write operation:**

W1      **choose newwptr such that** newwptr != oldwptr **and**
            newwptr != trace[r] **for all** r; /* *oldwptr* initialized to ⊥*/

W2      **write** BUF[newwptr];

W3      wsync := swap(**&SYNC**, 0 | newwptr); /* Clears all reading bits */

W4      oldwptr := newwptr;
         usedwptr := wsync **&** PTRFIELD;

W5      **for each reader** r

W6          **if** (wsync **&** (1 << (r + PTRFIELDLEN)))

W7              trace[r] := usedwptr;

# Atomic MRSW Register – ARC [Ianni]

---

**Algorithm 1.** Register Initialization

1: **procedure** INIT($content$, $size$)
2:     **for all** $slot \in [0, N+1]$ **do**
3:       $register[slot].size \leftarrow 0$
4:       $register[slot].r\_start \leftarrow 0$
5:       $register[slot].r\_end \leftarrow 0$
6:     MEMCOPY $register[0].content, content, size$
7:     $register[0].size \leftarrow size$
8:     $current \leftarrow N$

---

**Algorithm 2.** The Atomic Register Read Operation

1: **procedure** READ
2:     $index \leftarrow current \gg 32$       ▷ **R1**
3:     **if** $last\_index = index$ **then**
4:       $entry \leftarrow register[last\_index]$
5:       **return** $\langle entry.content, entry.size \rangle$   ▷ **R2**
6:     ATOMICINC($register[last\_index].r\_end$)   ▷ **R3**
7:     $tmp\_curr \leftarrow$ AtomicAddAndFetch $(current, 1)$  ▷ **R4**
8:     $last\_index \leftarrow tmp\_curr \gg 32$     ▷ **R5**
9:     $entry \leftarrow register[last\_index]$
10:    **return** $\langle entry.content, entry.size \rangle$

---

**Algorithm 3.** The Atomic Register Write Operation

1: **procedure** WRITE($content$, $size$)
2:     pick $slot$ such that $slot \neq last\_slot \wedge$
      $register[slot].r\_start = register[slot].r\_end$   ▷ **W1**
3:     MEMCOPY($register[slot].content, content, size$)
4:     $register[slot].size \leftarrow size$
5:     $register[slot].r\_start \leftarrow 0$
6:     $register[slot].r\_end \leftarrow 0$
7:     $old\_curr \leftarrow$ ATOMICEXCHANGE $(current, slot \ll 32)$  ▷ **W2**
8:     $old\_slot \leftarrow old\_curr \gg 32$
9:     $register[old\_slot].r\_start \leftarrow old\_curr \,\&\, (2^{32} - 1)$  ▷ **W3**
10:    $last\_slot \leftarrow slot$