

Lecture #09

Index Concurrency Control



ADMINISTRIVIA

Project #1 is due October 1st @ 11:59pm

→ Special Office Hours: **Saturday Sept 30th @ 3pm-5pm**

Homework #2 is due Wed Oct 4th @ 11:59pm

Homework #3 is due Sun Oct 8th @ 11:59pm

Mid-Term Exam is Wednesday Oct 11th

→ During regular class time from 2:00-3:20pm

→ Please contact us if you need accommodations.

→ More details next week...

OBSERVATION

We (mostly) assumed all the data structures that we have discussed so far are single-threaded.

But a DBMS needs to allow multiple threads to safely access data structures to take advantage of additional CPU cores and hide disk I/O stalls.



They Don't Do This!

VOLTDB



redis

H-Store

CONCURRENCY CONTROL

A concurrency control protocol is the method that the DBMS uses to ensure “correct” results for concurrent operations on a shared object.

A protocol's correctness criteria can vary:

- **Logical Correctness:** Can a thread see the data that it is supposed to see?
- **Physical Correctness:** Is the internal representation of the object sound?

TODAY'S AGENDA

Latches Overview

Hash Table Latching

B+Tree Latching

Leaf Node Scans

LOCKS VS. LATCHES

Locks (Transactions)

- Protect the database's logical contents from other transactions.
- Held for transaction's duration.
- Need to be able to rollback changes.

Latches (Workers)

- Protect the critical sections of the DBMS's internal data structure from other workers (e.g., threads).
- Held for operation duration.
- Do not need to be able to rollback changes.

LOCKS VS. LATCHES

Lecture #15



Locks

Separate...	Transactions
Protect...	Database Contents
During...	Entire Transactions
Modes...	Shared, Exclusive, Update, Intention
Deadlock	Detection & Resolution
...by...	Waits-for, Timeout, Aborts
Kept in...	Lock Manager

Latches

Workers (threads, processes)
In-Memory Data Structures
Critical Sections
Read, Write
Avoidance
Coding Discipline
Protected Data Structure

LATCH MODES

Read Mode

- Multiple threads can read the same object at the same time.
- A thread can acquire the read latch if another thread has it in read mode.

Write Mode

- Only one thread can access the object.
- A thread cannot acquire a write latch if another thread has it in any mode.

Compatibility Matrix

	Read	Write
Read	✓	X
Write	X	X

LATCH IMPLEMENTATION GOALS

Small memory footprint.

Fast execution path when no contention.

Deschedule thread when it has been waiting for too long to avoid burning cycles.

Each latch should not have to implement their own queue to track waiting threads.

LATCH IM

Small memory

Fast execution

Deschedule the
too long to avoid

Each latch should
own queue to

By: **Linus Torvalds** (torvalds.delete@this.linux-foundation.org), January 3, 2020 6:05 pm
 Beastian (no.email.delete@this.aol.com) on January 3, 2020 11:46 am wrote:
 > I'm usually on the other side of these primitives when I write code as a consumer of them,
 > but it's very interesting to read about the nuances related to their implementations:

Room: [Moderated Discussions](#)

The whole post seems to be just wrong, and is measuring something completely different than what the author thinks and claims it is measuring.

First off, spinlocks can only be used if you actually know you're not being scheduled while using them. But the blog post author seems to be implementing his own spinlocks in user space with no regard for whether the lock user might be scheduled or not. And the code used for the claimed "lock not held" timing is complete garbage.

It basically reads the time before releasing the lock, and then it reads it after acquiring the lock again, and claims that the time difference is the time when no lock was held. Which is just inane and pointless and completely wrong.

That's pure garbage. What happens is that

- since you're spinning, you're using CPU time
- at a random time, the scheduler will schedule you out
- that random time might be just after you read the "current time", but before you actually released the spinlock.

So now you still hold the lock, but you got scheduled away from the CPU, because you had used up your time slice. The "current time" you read is basically now stale, and has nothing to do with the (future) time when you are *actually* going to release the lock.

Somebody else comes in and wants that "spinlock", and that somebody will now spin for a long while, since nobody is releasing it - it's still held by that other thread entirely that was just scheduled out. At some point, the scheduler says "ok, now you've used your time slice", and schedules the original thread, and *now* the lock is actually released. Then another thread comes in, gets the lock again, and then it looks at the time and says "oh, a long time passed without the lock being held at all".

And notice how the above is the *good* scenario. If you have more threads than CPU's (maybe because of other processes unrelated to your own test load), maybe the next thread that gets scheduled isn't the one that is going to release the lock. No, that one already got its timeslice, so the next thread scheduled might be *another* thread that wants that lock that is still being held by the thread that isn't even running right now!

So the code in question is pure garbage. You can't do spinlocks like that. Or rather, you very much can do them like that, and when you do that you are measuring random latencies and getting nonsensical values, because what you are measuring is "I have a lot of busywork, where all the processes are CPU-bound, and I'm measuring random points of how long the scheduler kept the process in place".

And then you write a blog-post blaming others, not understanding that it's your incorrect code that is garbage, and is giving random garbage values.

LATCH IM

Small memory

Fast execution

Deschedule th

I repeat: **do not use spinlocks in user space, unless you actually know what you're doing.** And be aware that the likelihood that you know what you are doing is basically nil.

By: **Linus Torvalds** (torvalds.delete@this.linux-foundation.org), January 3, 2020 6:05 pm

Beastian (no.email.delete@this.aol.com) on January 3, 2020 11:46 am wrote:
> I'm usually on the other side of these primitives when I write code as a consumer of them,
> but it's very interesting to read about the nuances related to their implementations:

Room: [Moderated Discussions](#)

The whole post seems to be just wrong, and is measuring something completely different than what the author thinks and claims it is measuring.

First off, spinlocks can only be used if you actually know you're not being scheduled while using them. But the blog post author seems to be implementing his own spinlocks in user space with no regard for whether the lock user might be scheduled or not. And the code used for the claimed "lock not held" timing is complete garbage.

It basically reads the time before releasing the lock, and then it reads it after acquiring the lock again, and claims that the time difference is the time when no lock was held. Which is just inane and pointless and completely wrong.

That's pure garbage. What happens is that

- (a) since you're spinning, you're using CPU time
- (b) at a random time, the scheduler will schedule you out
- (c) that random time might be just after you...

...of the processes are CPU-bound, and getting nonsensical values, because what you are measuring is "I have a lot of busywork, and then you write a blog-post blaming others, not understanding that it's your incorrect code that is garbage, and is giving random garbage values."

LATCH IMPLEMENTATIONS

Test-and-Set Spinlock

Blocking OS Mutex

Reader-Writer Locks

Advanced approaches:

→ Adaptive Spinlock ([Apple ParkingLot](#))

→ Queue-based Spinlock ([MCS Locks](#))

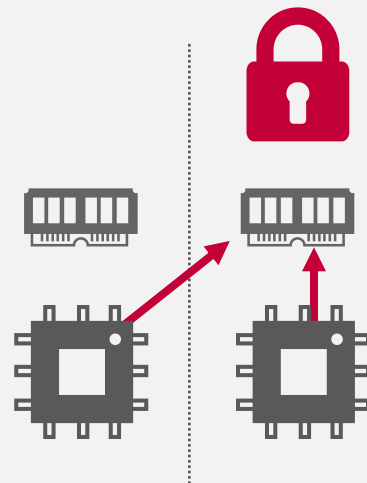
LATCH IMPLEMENTATIONS

Approach #1: Test-and-Set Spin Latch (TAS)

- Very efficient (single instruction to latch/unlatch)
- Non-scalable, not cache friendly, not OS friendly.
- Example: `std::atomic<T>`

std::atomic<bool>

```
std::atomic_flag latch;  
:  
while (latch.test_and_set(...)) {  
    // Retry? Yield? Abort?  
}
```

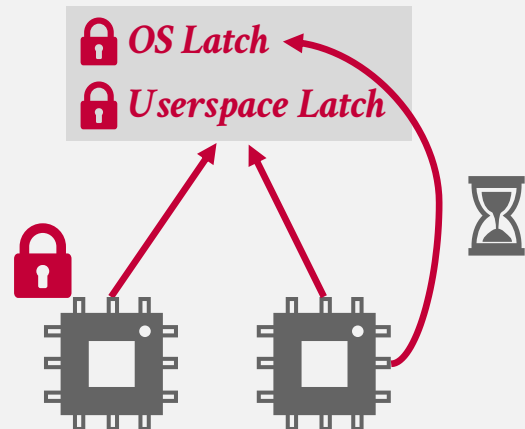


LATCH IMPLEMENTATIONS

Approach #2: Blocking OS Mutex

- Simple to use
- Non-scalable (about 25ns per lock/unlock invocation)
- Example: `std::mutex` → `pthread_mutex` → `futex`

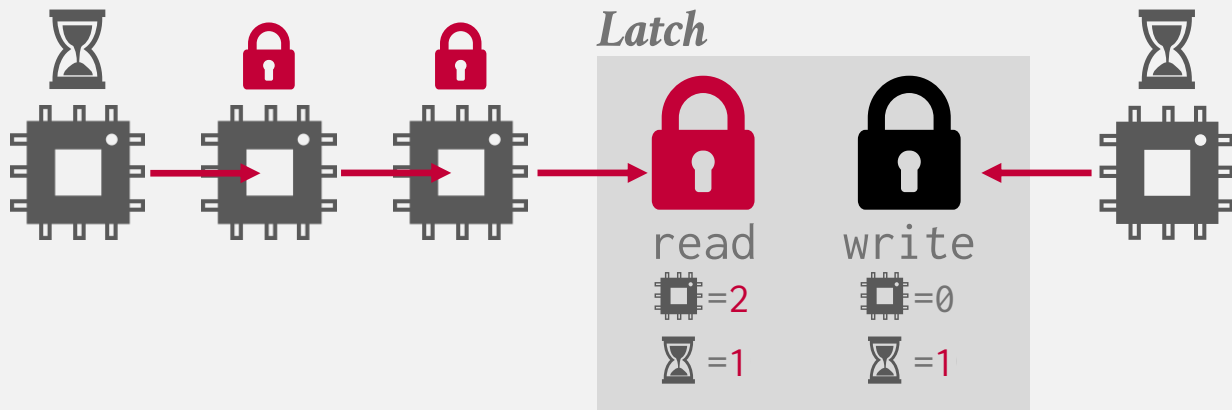
```
std::mutex m;  
:  
m.lock();  
// Do something special...  
m.unlock();
```



LATCH IMPLEMENTATIONS

Approach #3: Reader-Writer Latches

- Allows for concurrent readers. Must manage read/write queues to avoid starvation.
- Can be implemented on top of spinlocks.
- Example: `std::shared_mutex` → `pthread_rwlock`



HASH TABLE LATCHING

Easy to support concurrent access due to the limited ways threads access the data structure.

- All threads move in the same direction and only access a single page/slot at a time.
- Deadlocks are not possible.

To resize the table, take a global write latch on the entire table (e.g., in the header page).

HASH TABLE LATCHING

Approach #1: Page Latches

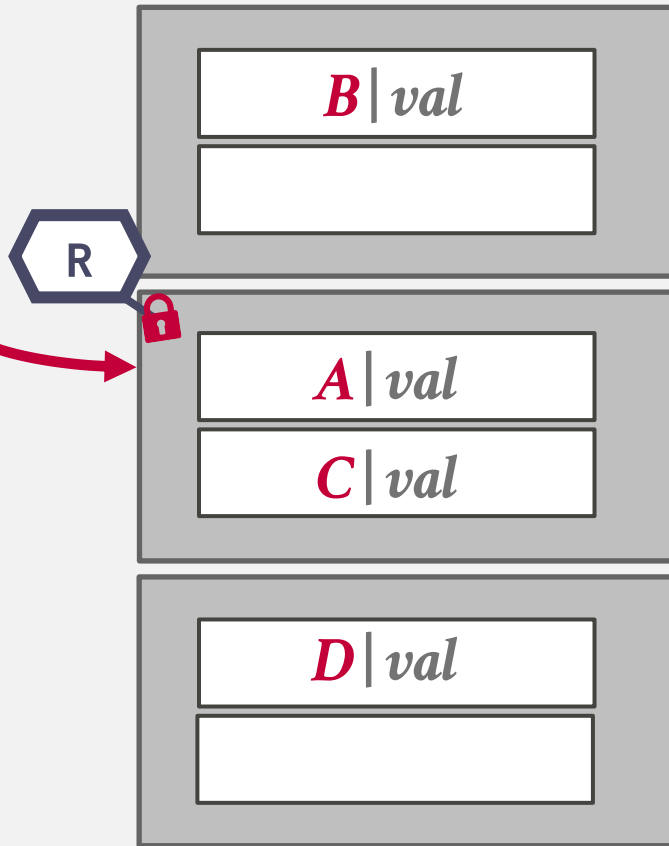
- Each page has its own reader-writer latch that protects its entire contents.
- Threads acquire either a read or write latch before they access a page.

Approach #2: Slot Latches

- Each slot has its own latch.
- Can use a single-mode latch to reduce meta-data and computational overhead.

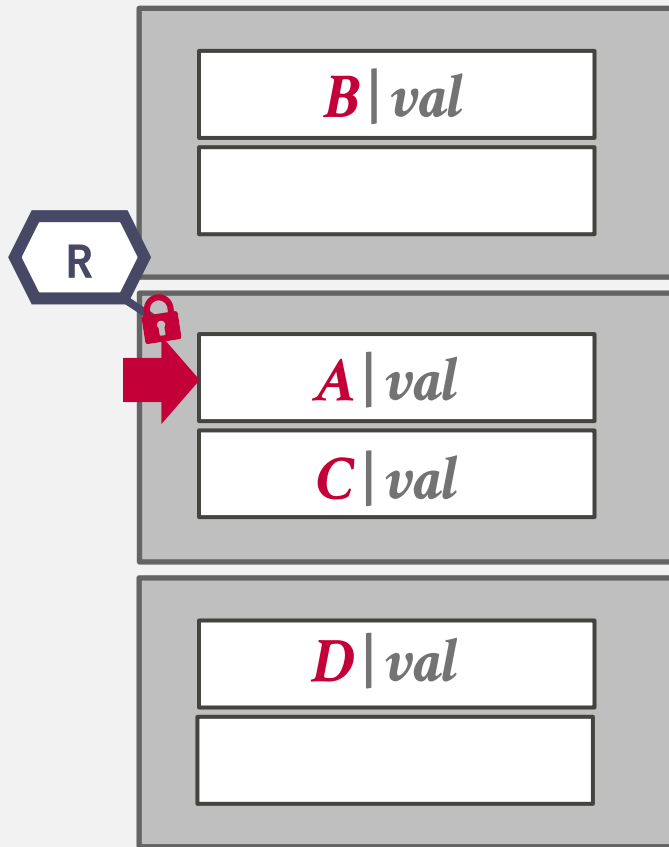
HASH TABLE - PAGE LATCHES

T_1 : Find D
 $hash(D)$



HASH TABLE - PAGE LATCHES

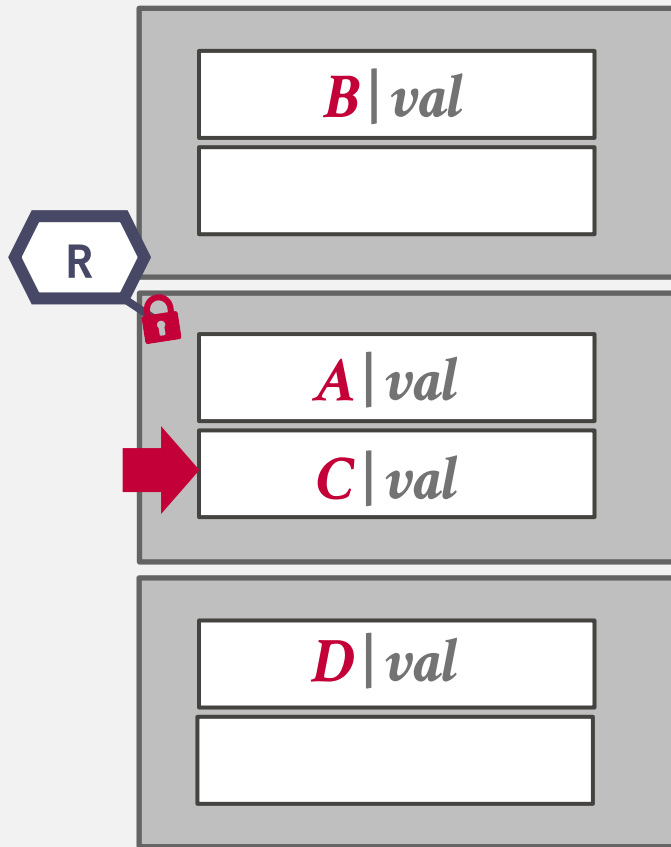
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

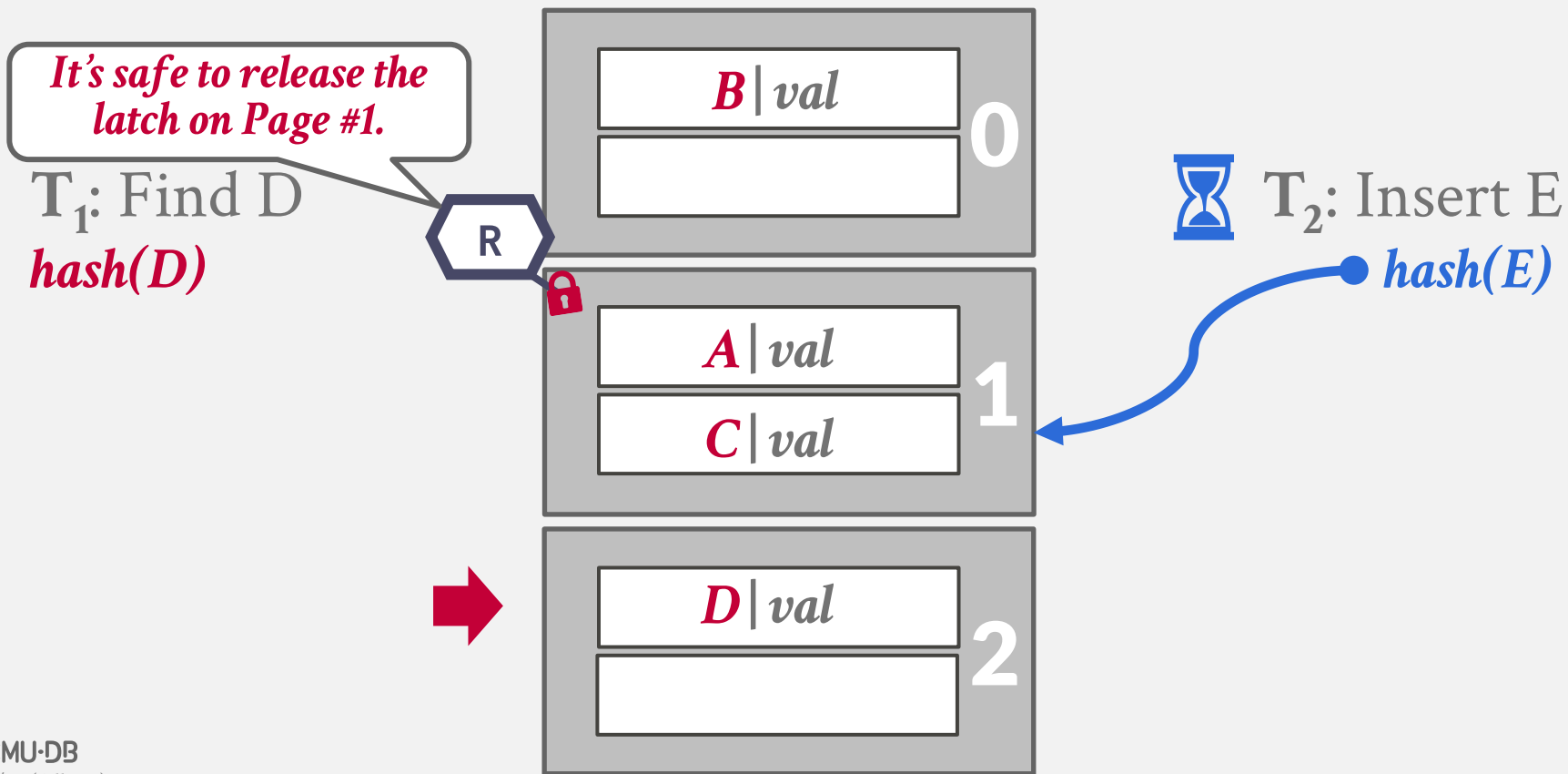
HASH TABLE - PAGE LATCHES

T_1 : Find D
hash(D)



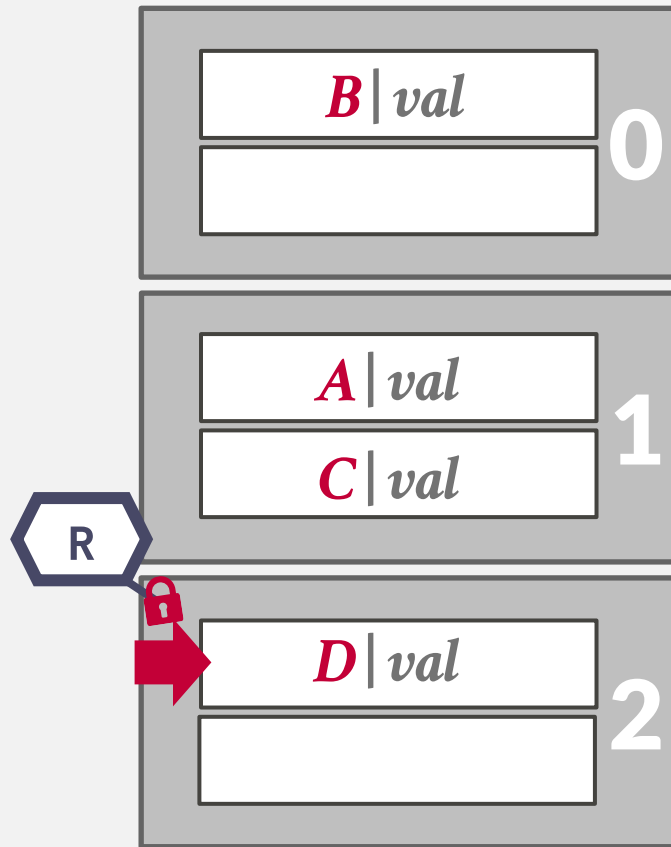
T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES



HASH TABLE - PAGE LATCHES

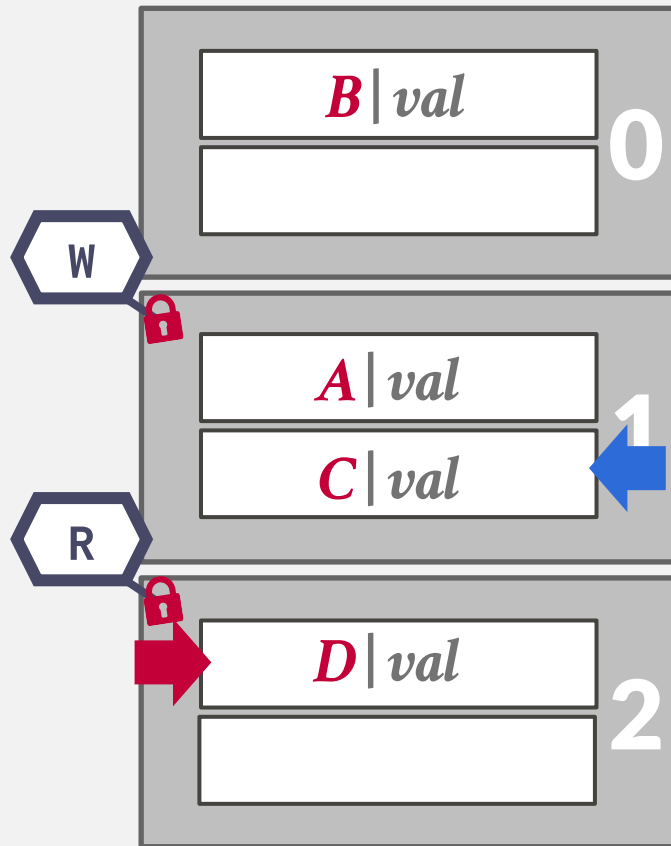
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES

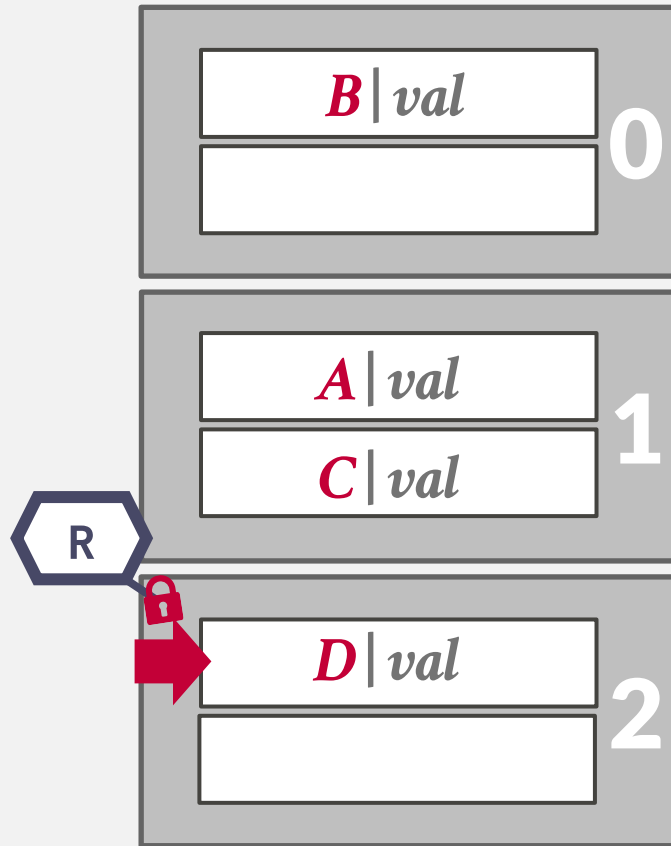
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES

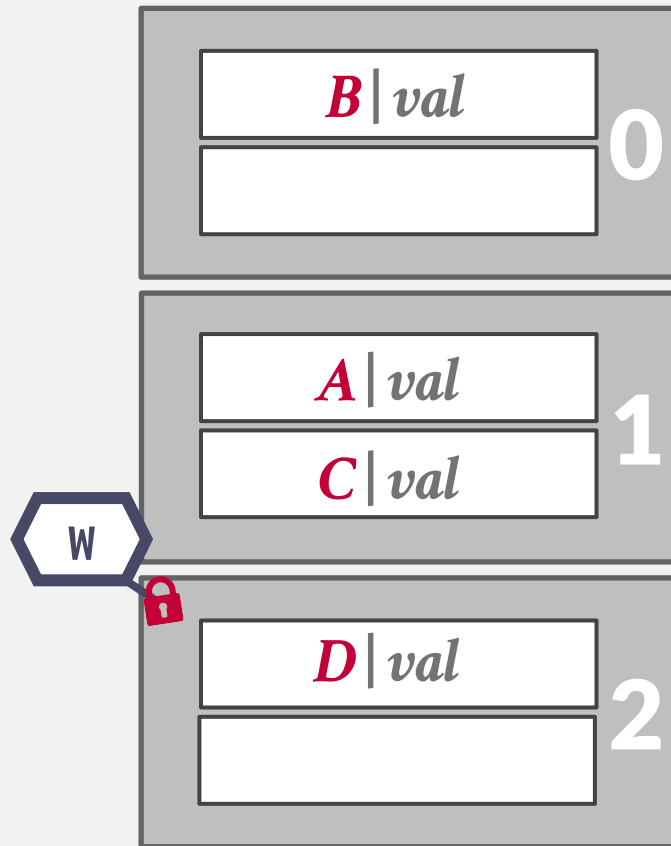
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES

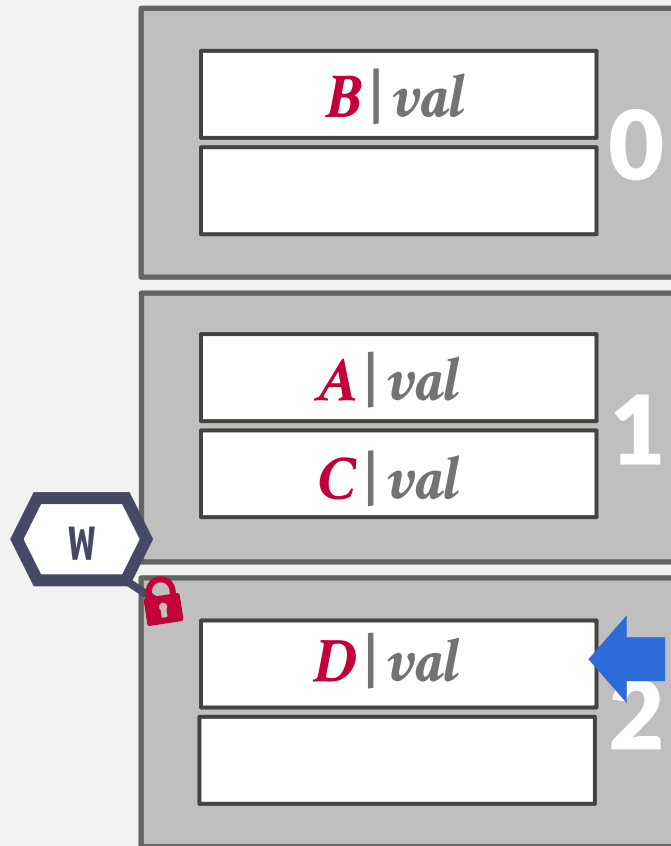
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES

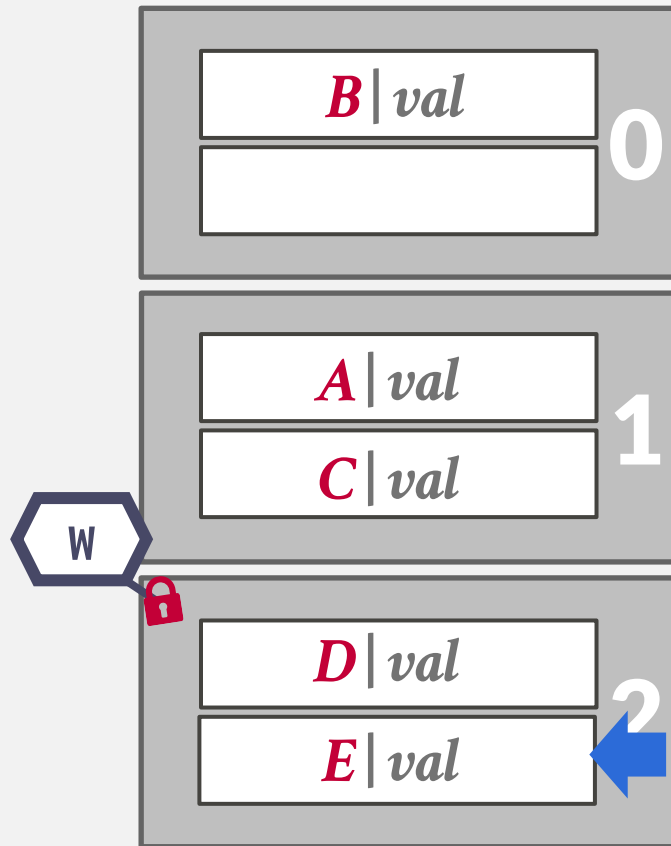
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - PAGE LATCHES

T_1 : Find D
hash(D)

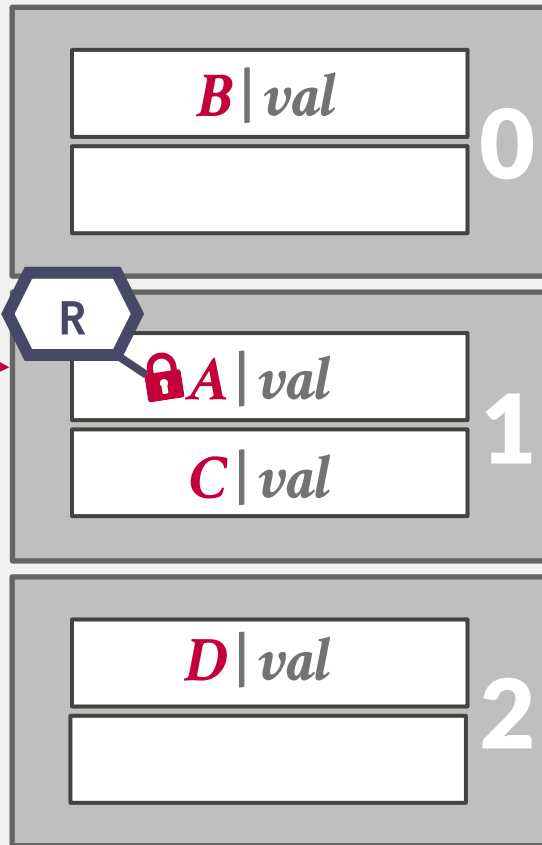


T_2 : Insert E
hash(E)

HASH TABLE - SLOT LATCHES

T_1 : Find D

$hash(D)$

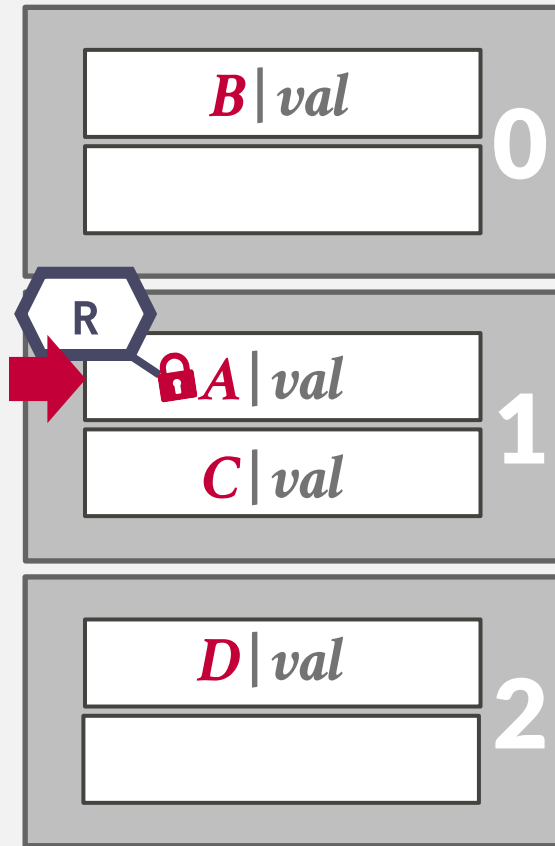


T_2 : Insert E

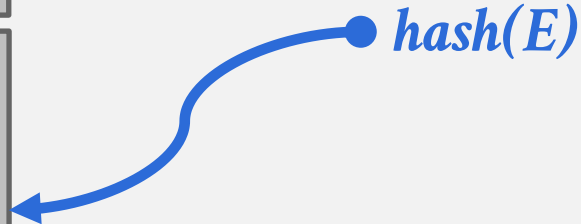
$hash(E)$

HASH TABLE - SLOT LATCHES

T_1 : Find D
hash(D)

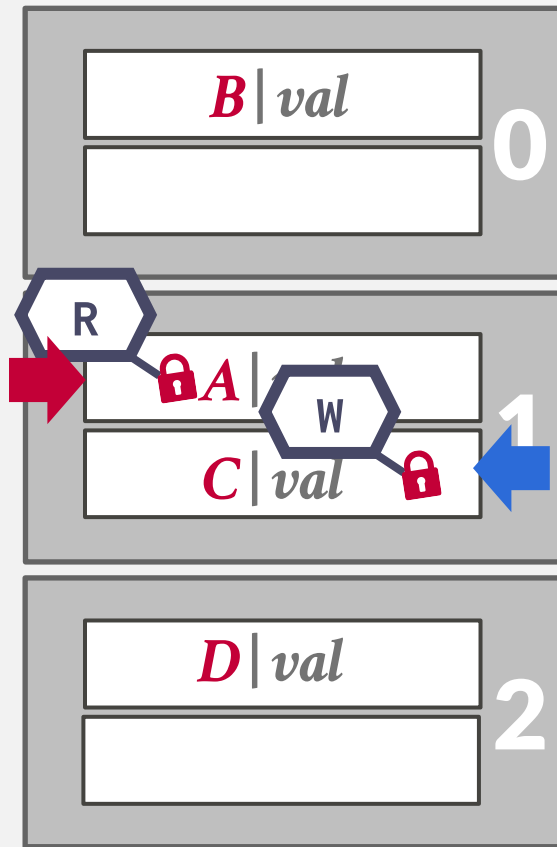


T_2 : Insert E
hash(E)



HASH TABLE - SLOT LATCHES

T_1 : Find D
hash(D)



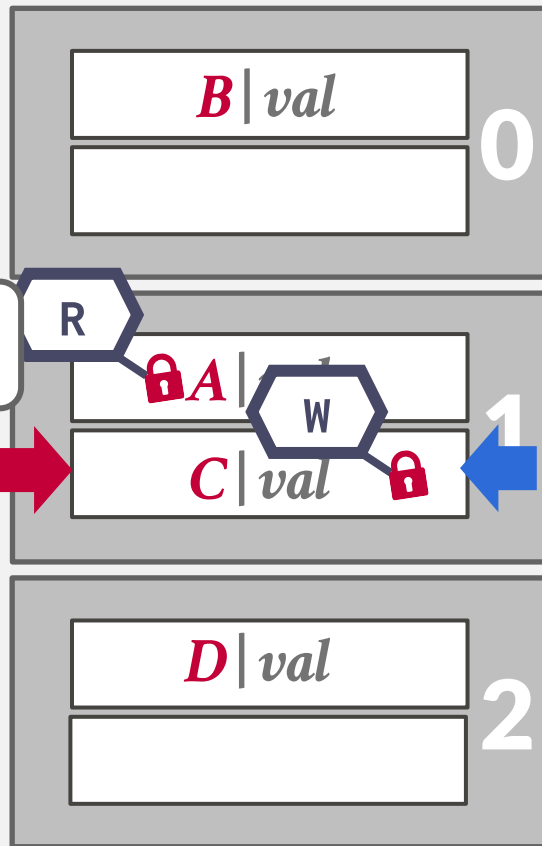
T_2 : Insert E
hash(E)

HASH TABLE - SLOT LATCHES

T_1 : Find D

$hash(D)$

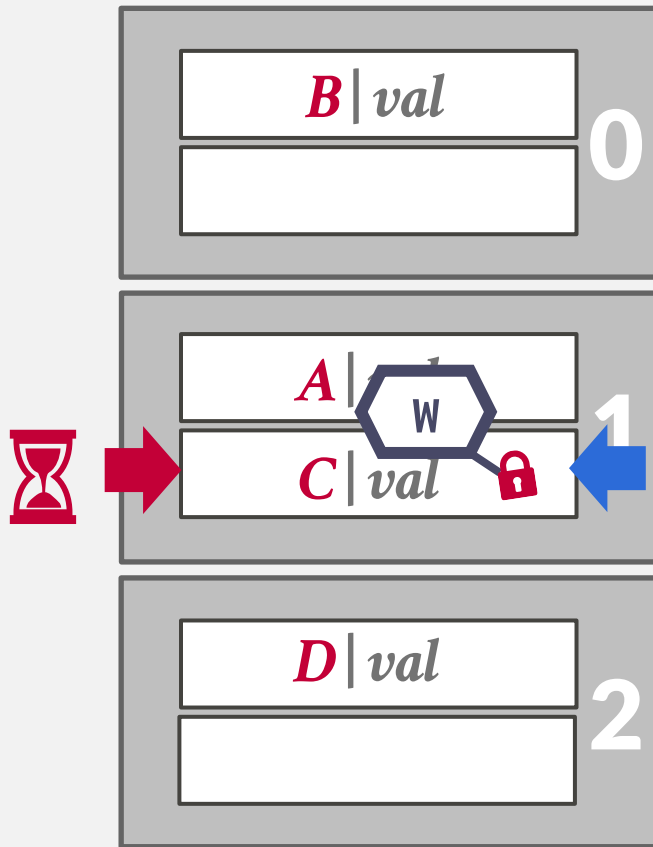
It's safe to release the latch on A



T_2 : Insert E
 $hash(E)$

HASH TABLE - SLOT LATCHES

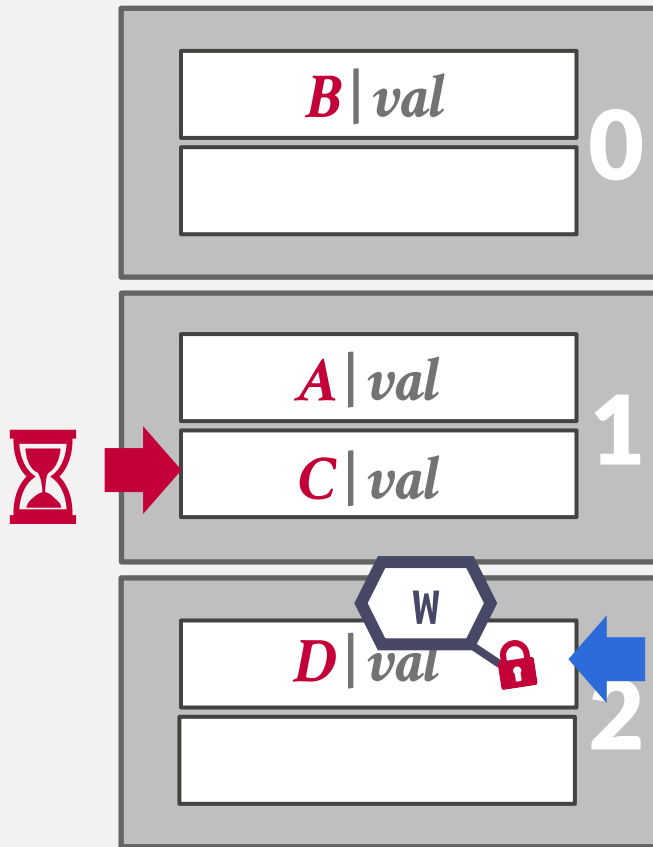
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - SLOT LATCHES

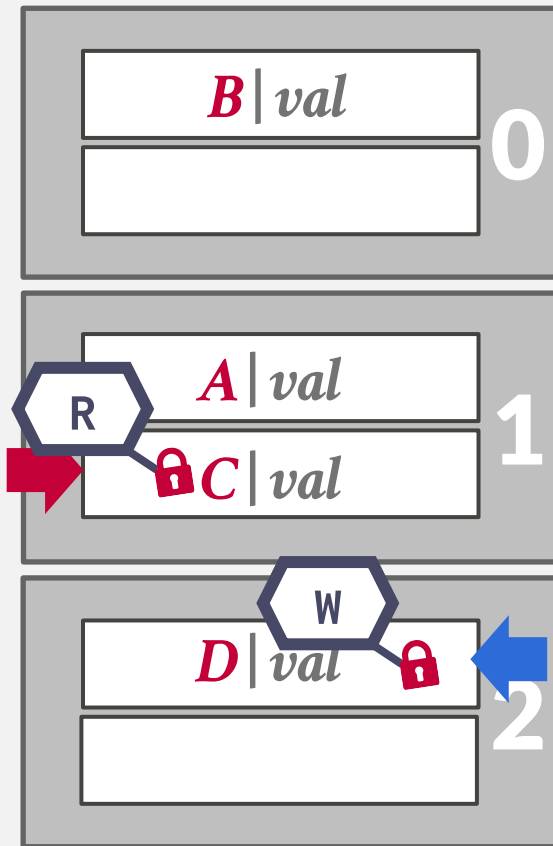
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

HASH TABLE - SLOT LATCHES

T_1 : Find D
hash(D)

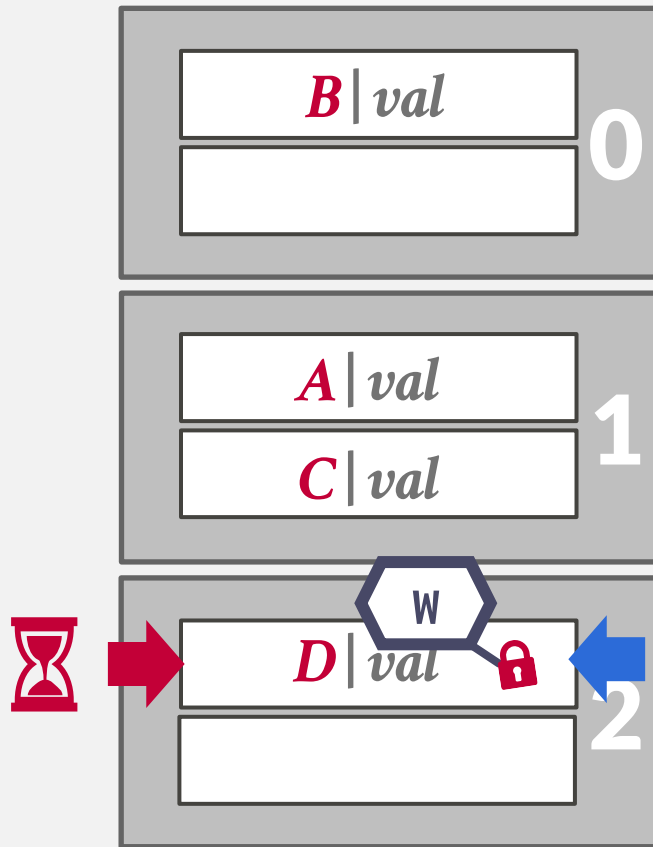


T_2 : Insert E
hash(E)

HASH TABLE - SLOT LATCHES

T_1 : Find D
hash(D)

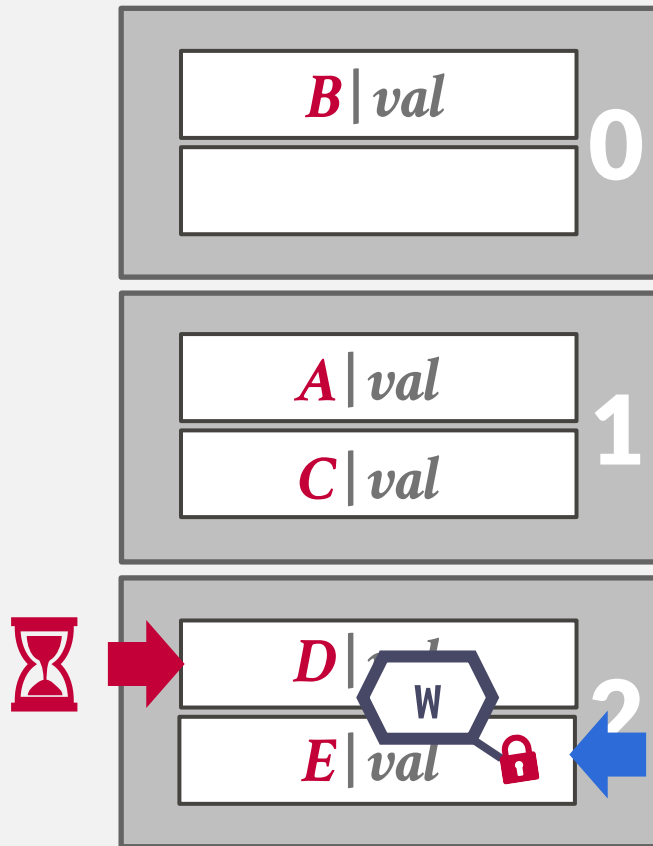
T_2 : Insert E
hash(E)



HASH TABLE - SLOT LATCHES

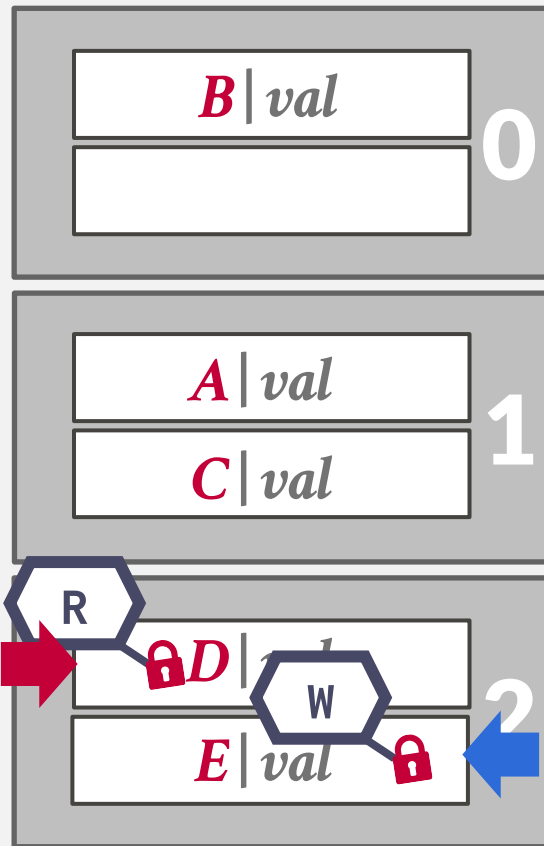
T_1 : Find D
hash(D)

T_2 : Insert E
hash(E)



HASH TABLE - SLOT LATCHES

T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

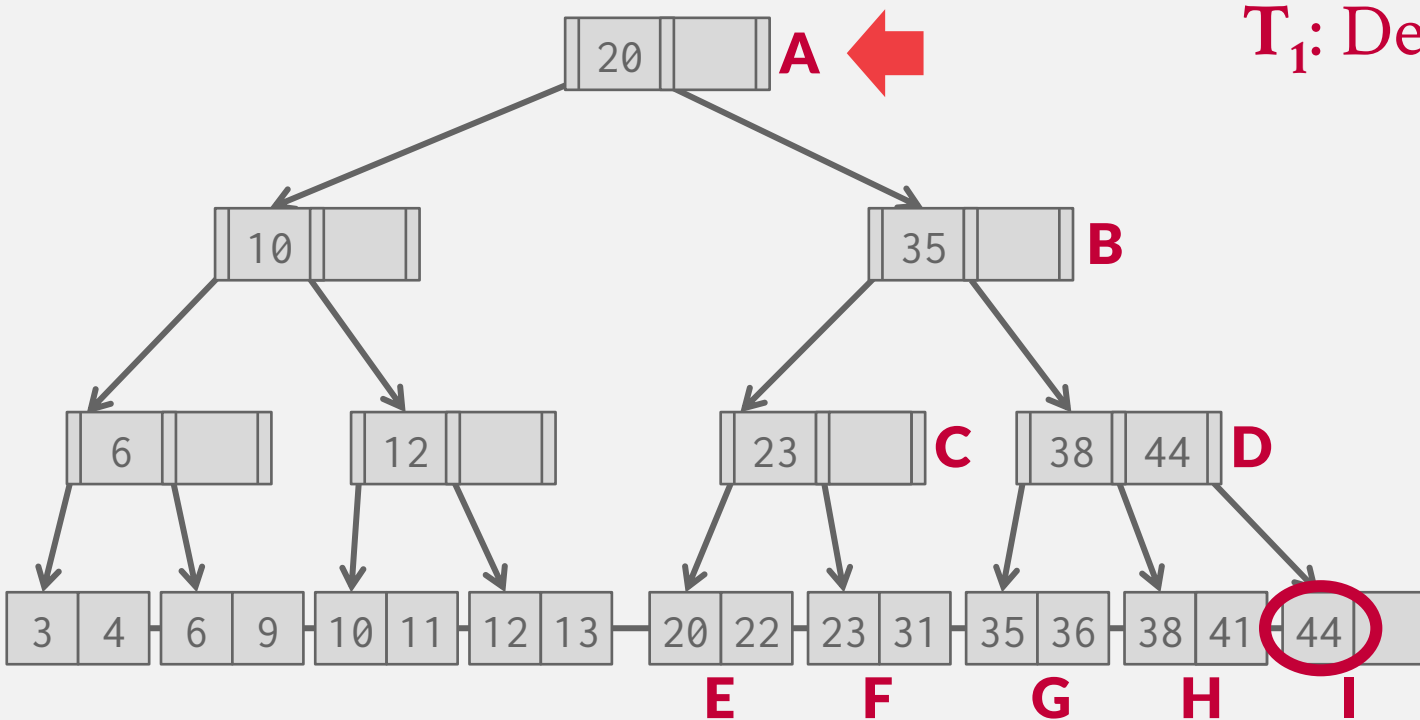
B+TREE CONCURRENCY CONTROL

We want to allow multiple threads to read and update a B+Tree at the same time.

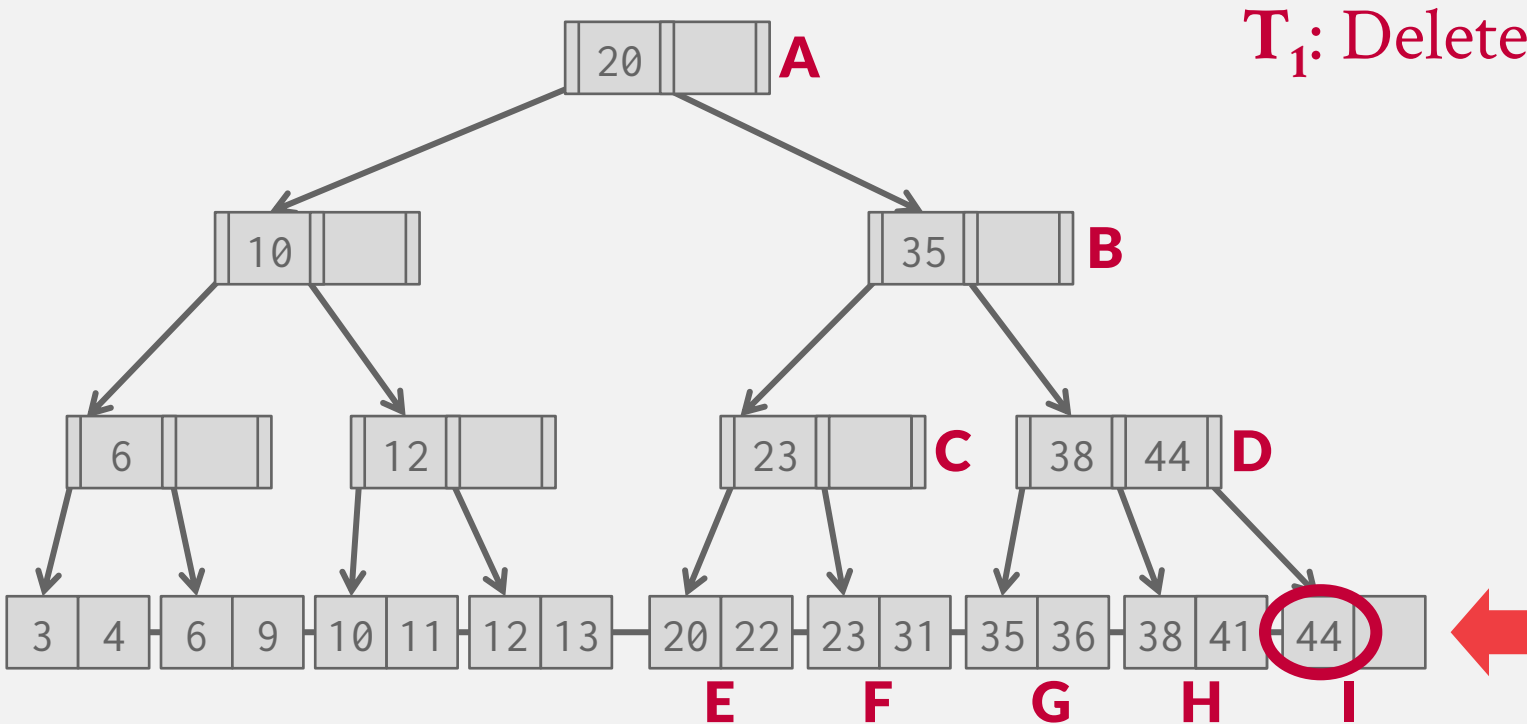
We need to protect against two types of problems:

- Threads trying to modify the contents of a node at the same time.
- One thread traversing the tree while another thread splits/merges nodes.

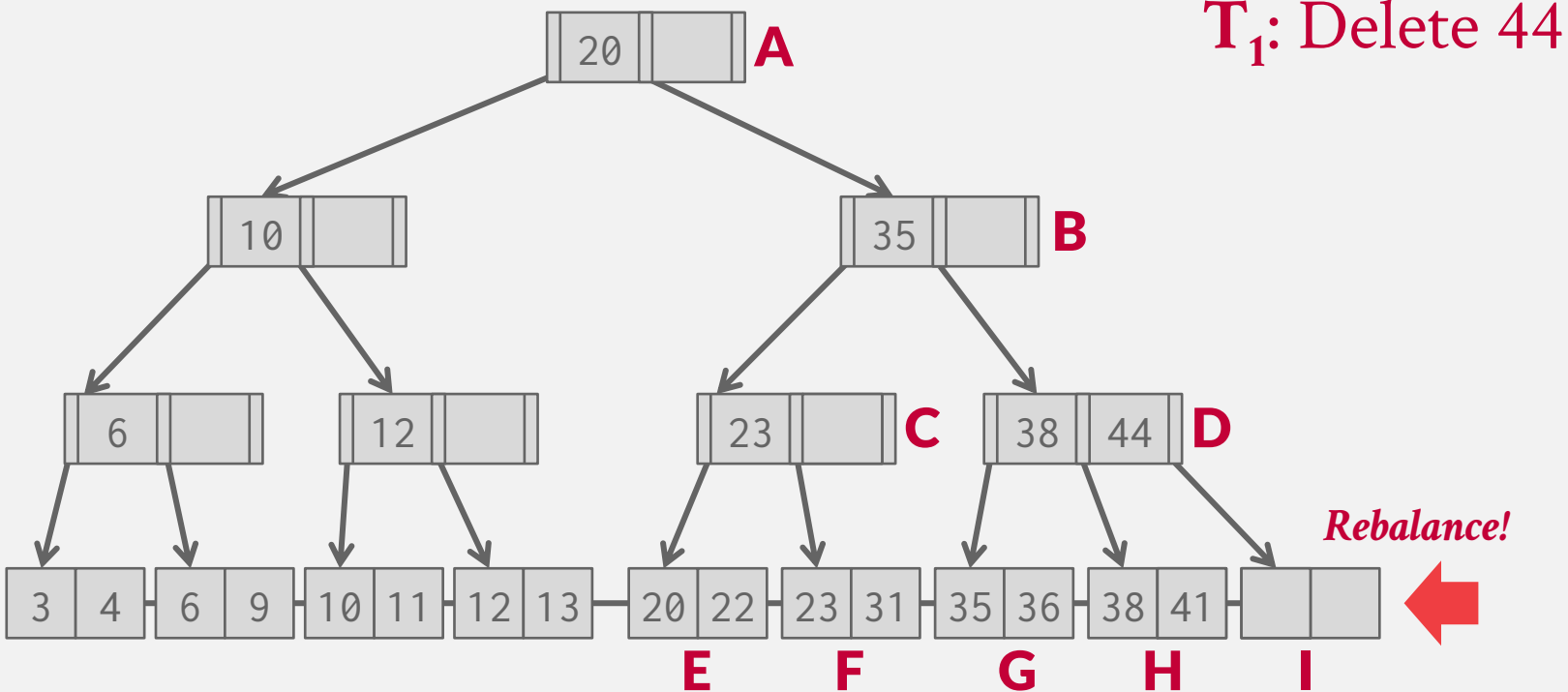
B+TREE MULTI-THREADED EXAMPLE



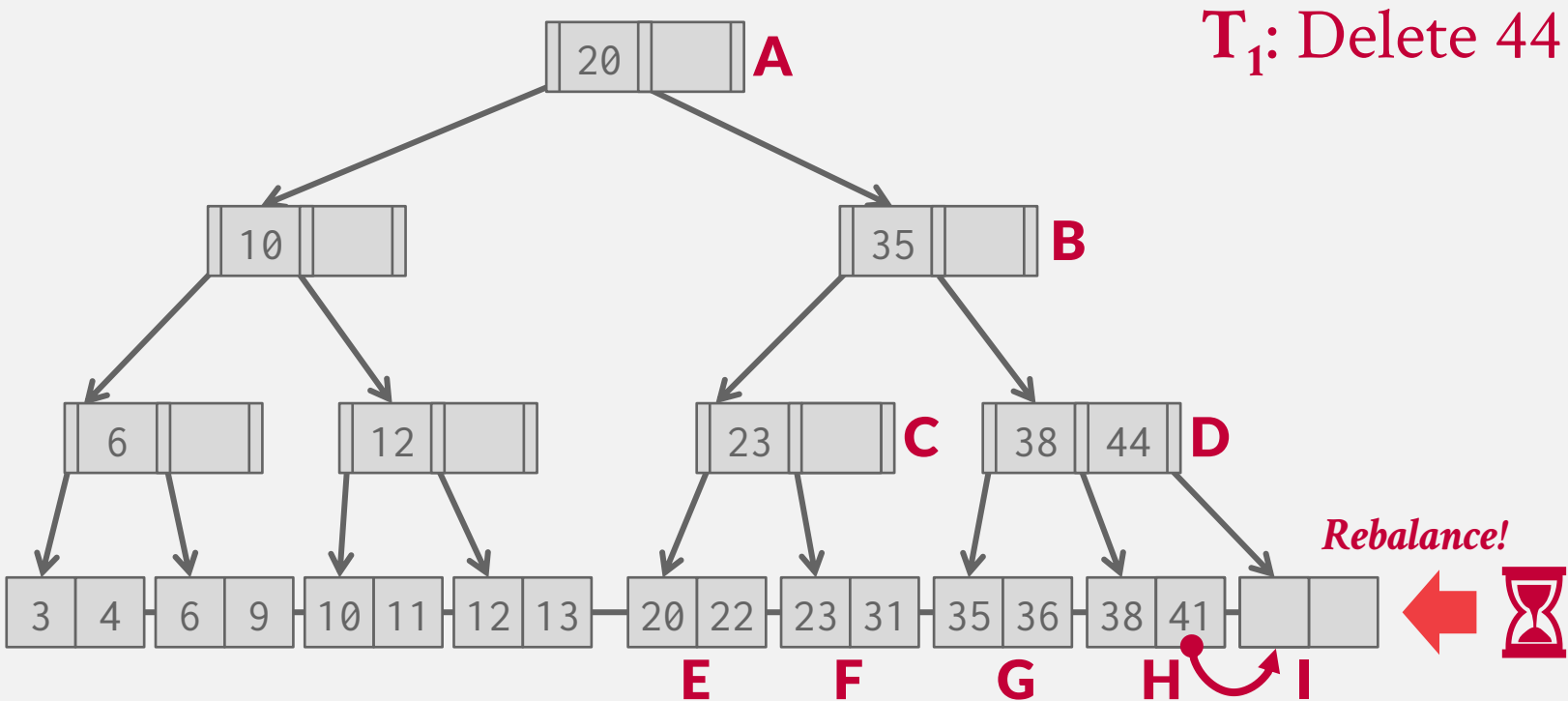
B+TREE MULTI-THREADED EXAMPLE



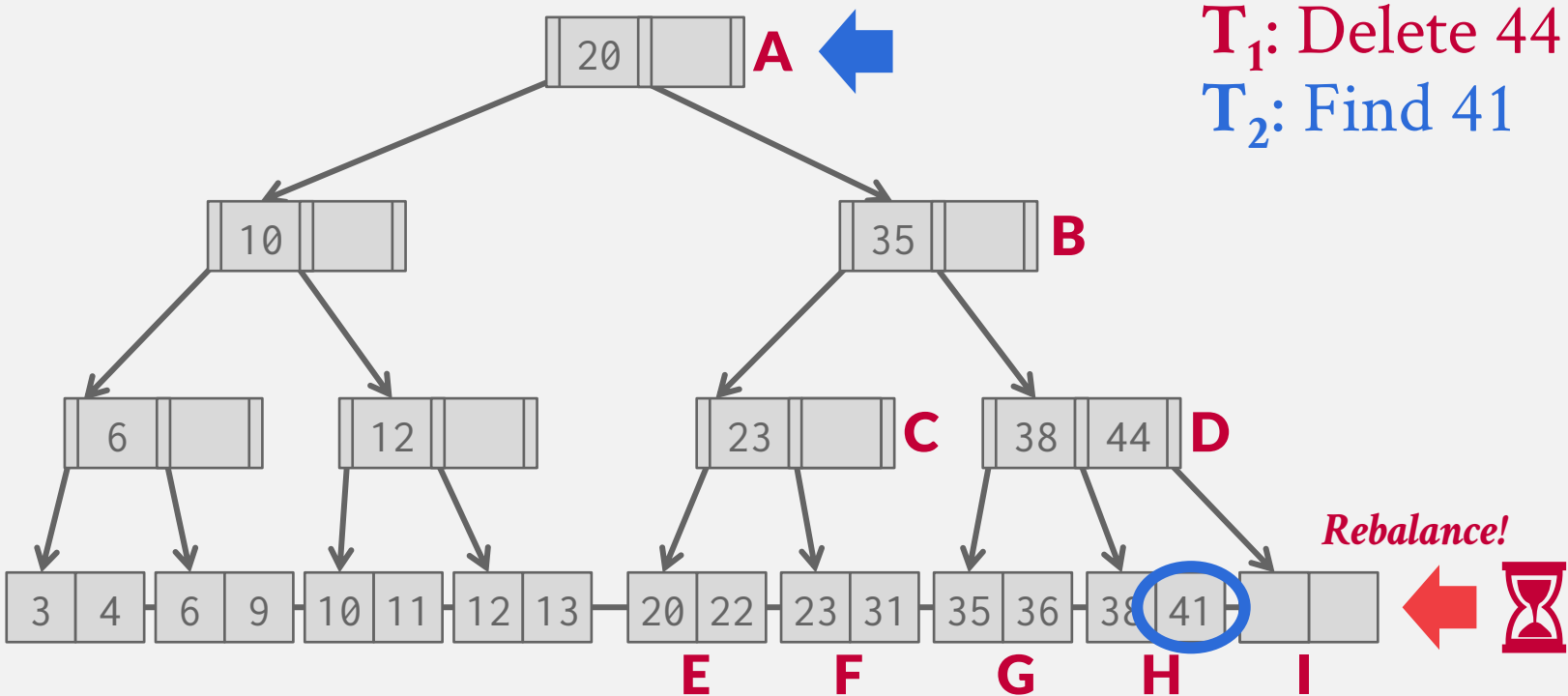
B+TREE MULTI-THREADED EXAMPLE



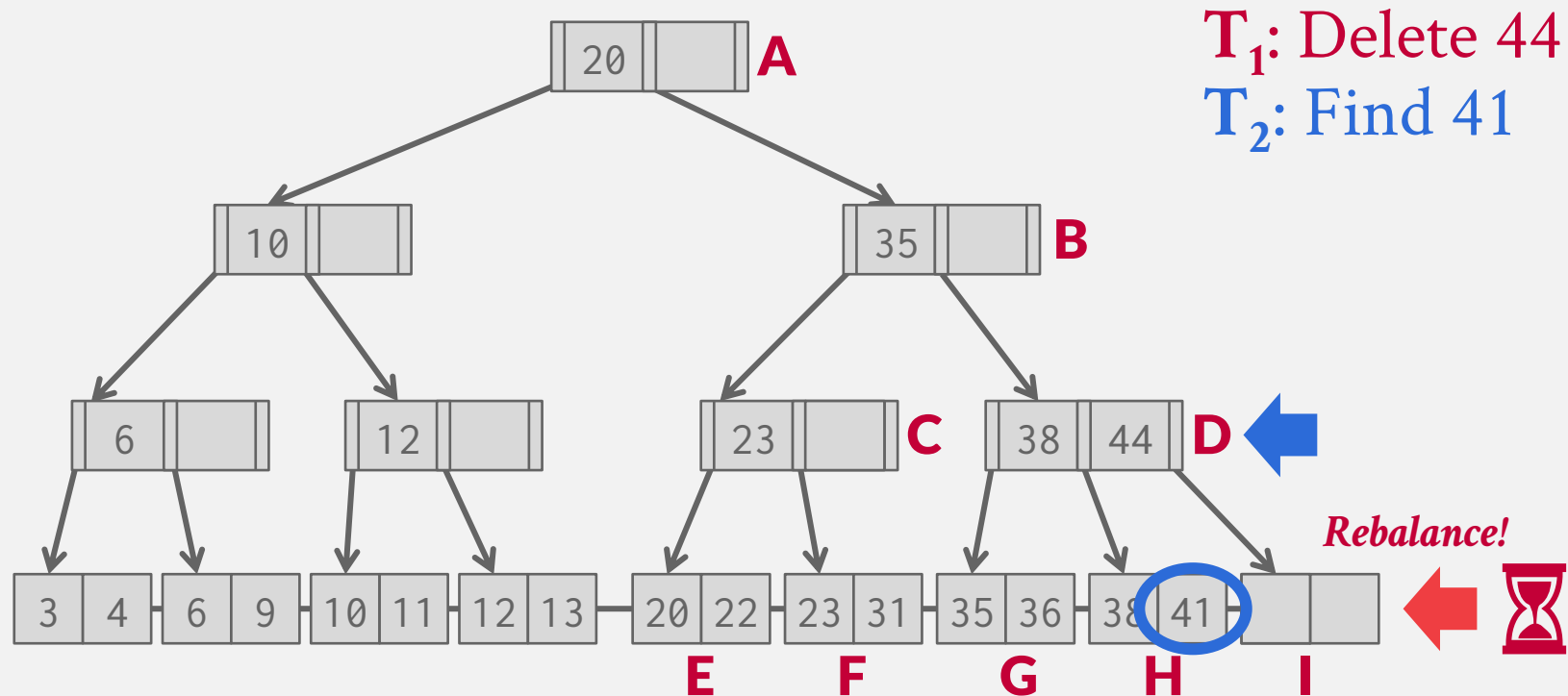
B+TREE MULTI-THREADED EXAMPLE



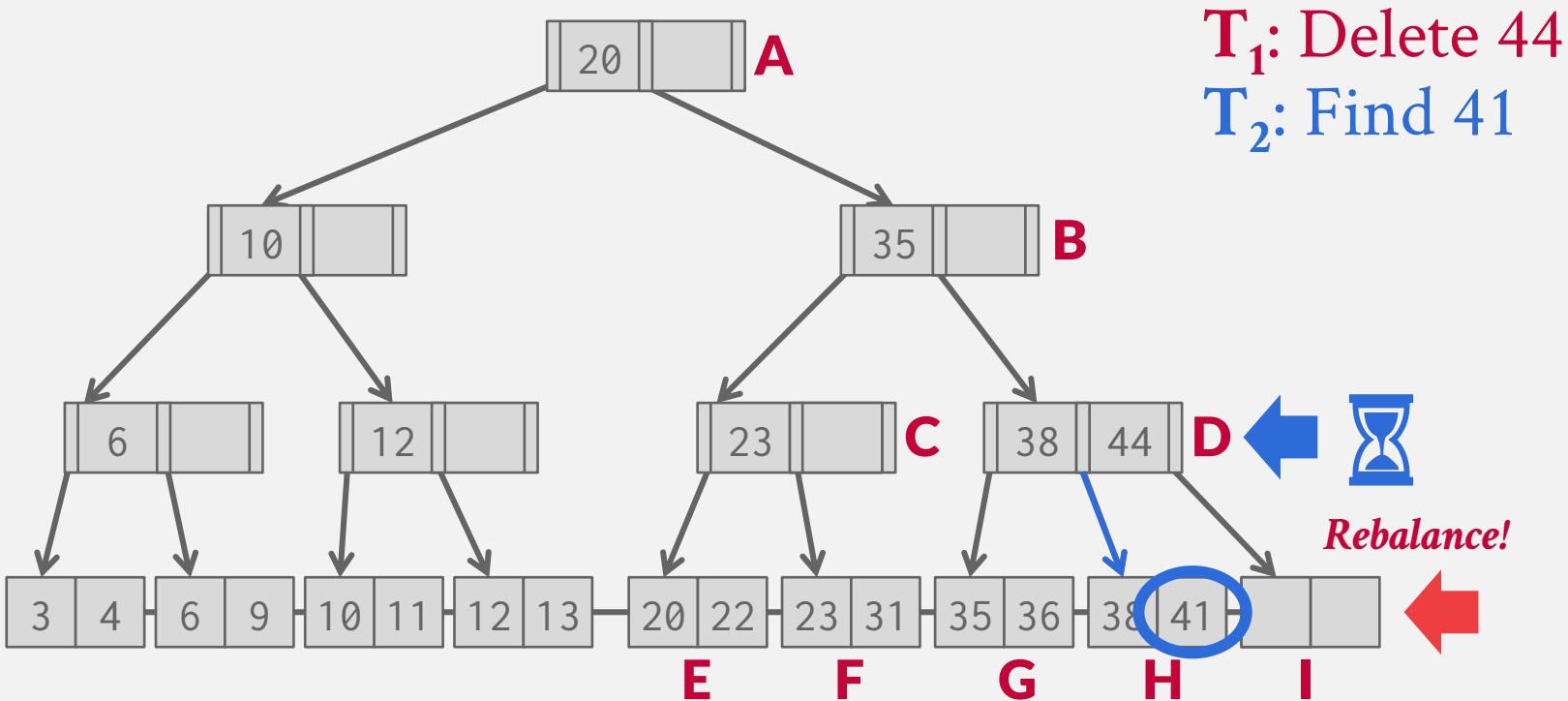
B+TREE MULTI-THREADED EXAMPLE



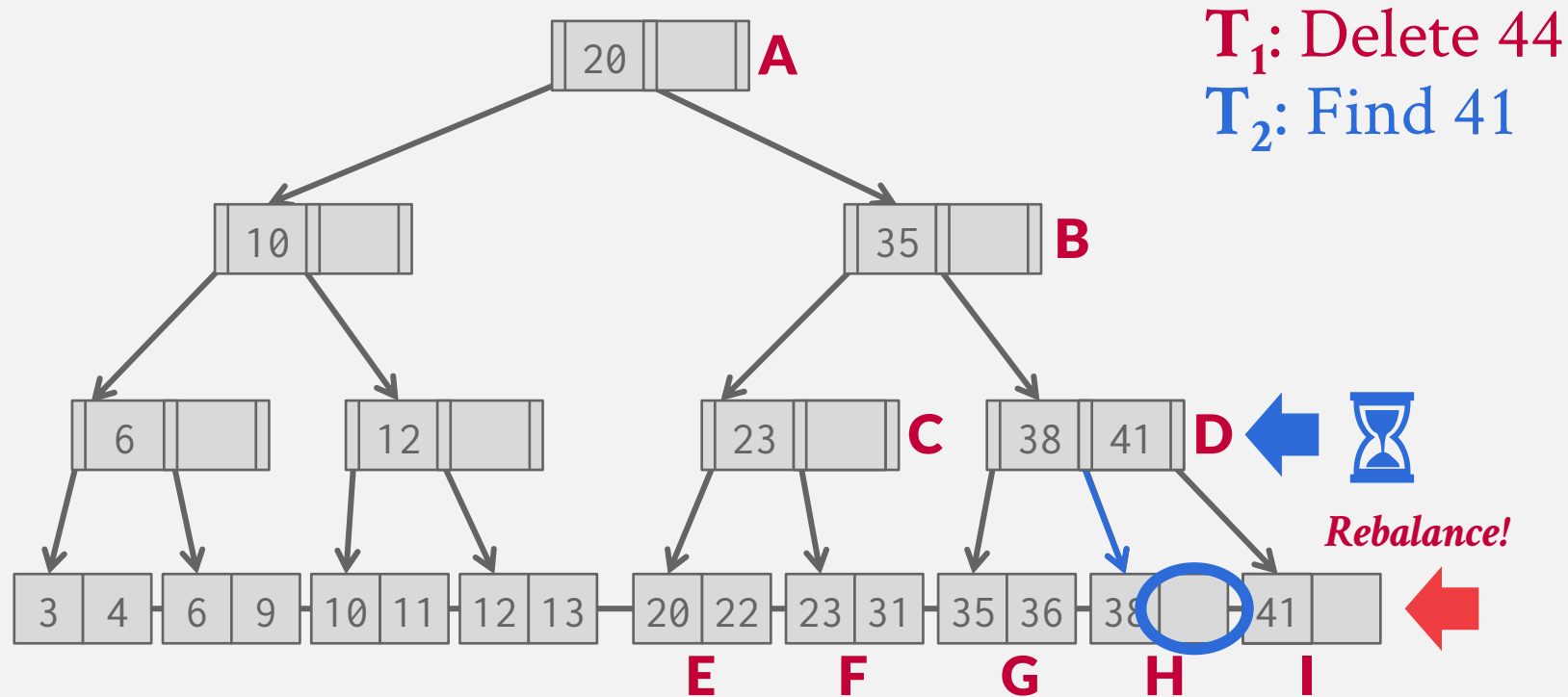
B+TREE MULTI-THREADED EXAMPLE



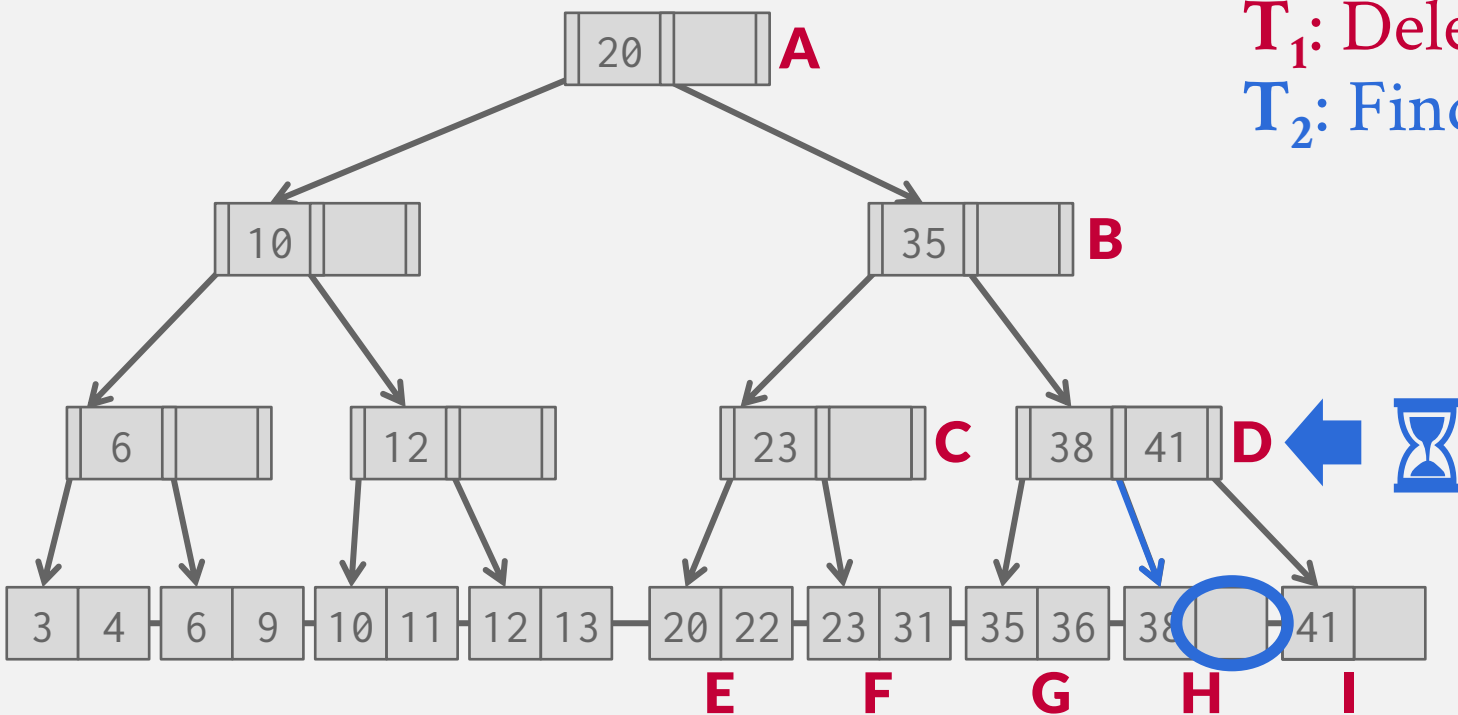
B+TREE MULTI-THREADED EXAMPLE



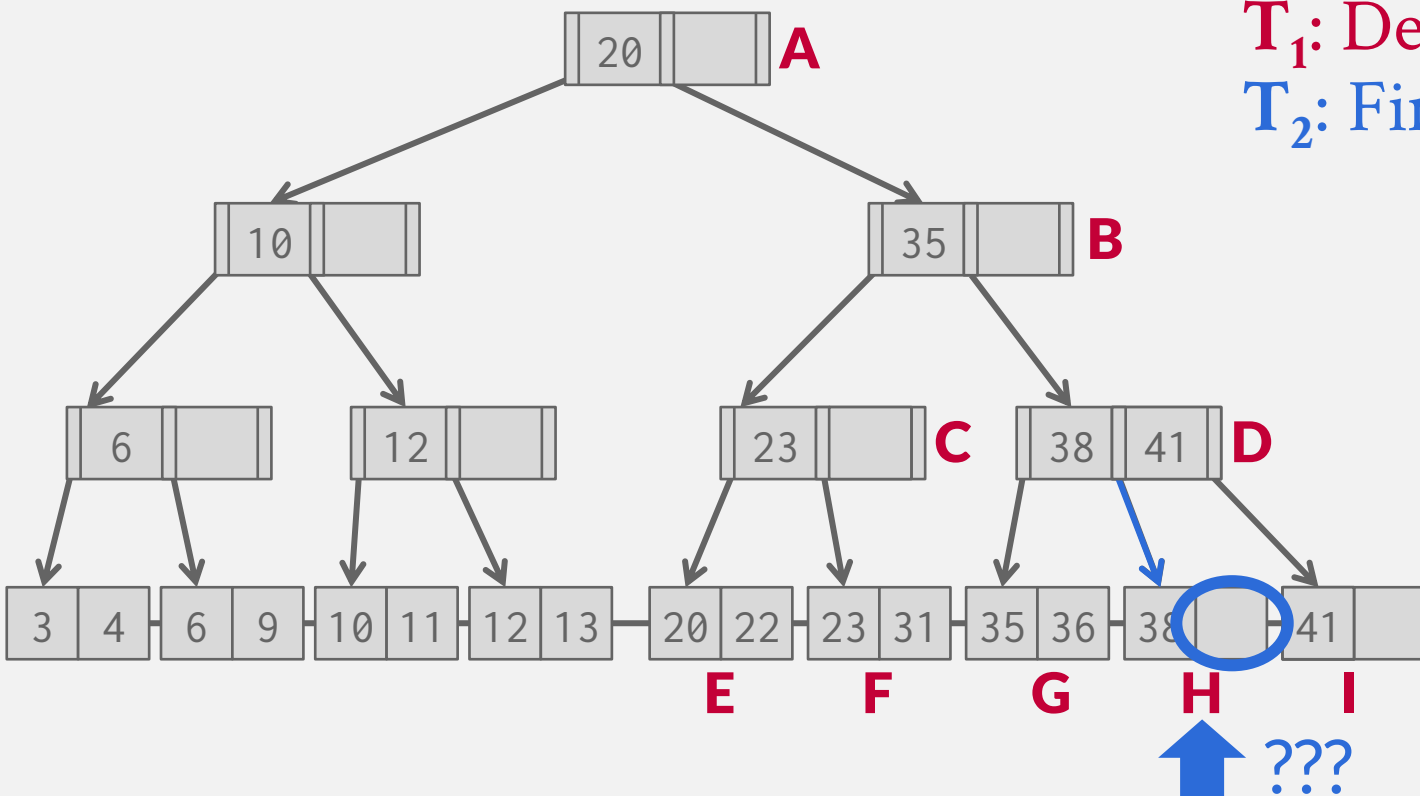
B+TREE MULTI-THREADED EXAMPLE



B+TREE MULTI-THREADED EXAMPLE



B+TREE MULTI-THREADED EXAMPLE



LATCH CRABBING/COUPLING

Protocol to allow multiple threads to access/modify B+Tree at the same time.

- Get latch for parent
- Get latch for child
- Release latch for parent if “safe”

A **safe node** is one that will not split or merge when updated.

- Not full (on insertion)
- More than half-full (on deletion)

LATCH CRABBING/COUPLING

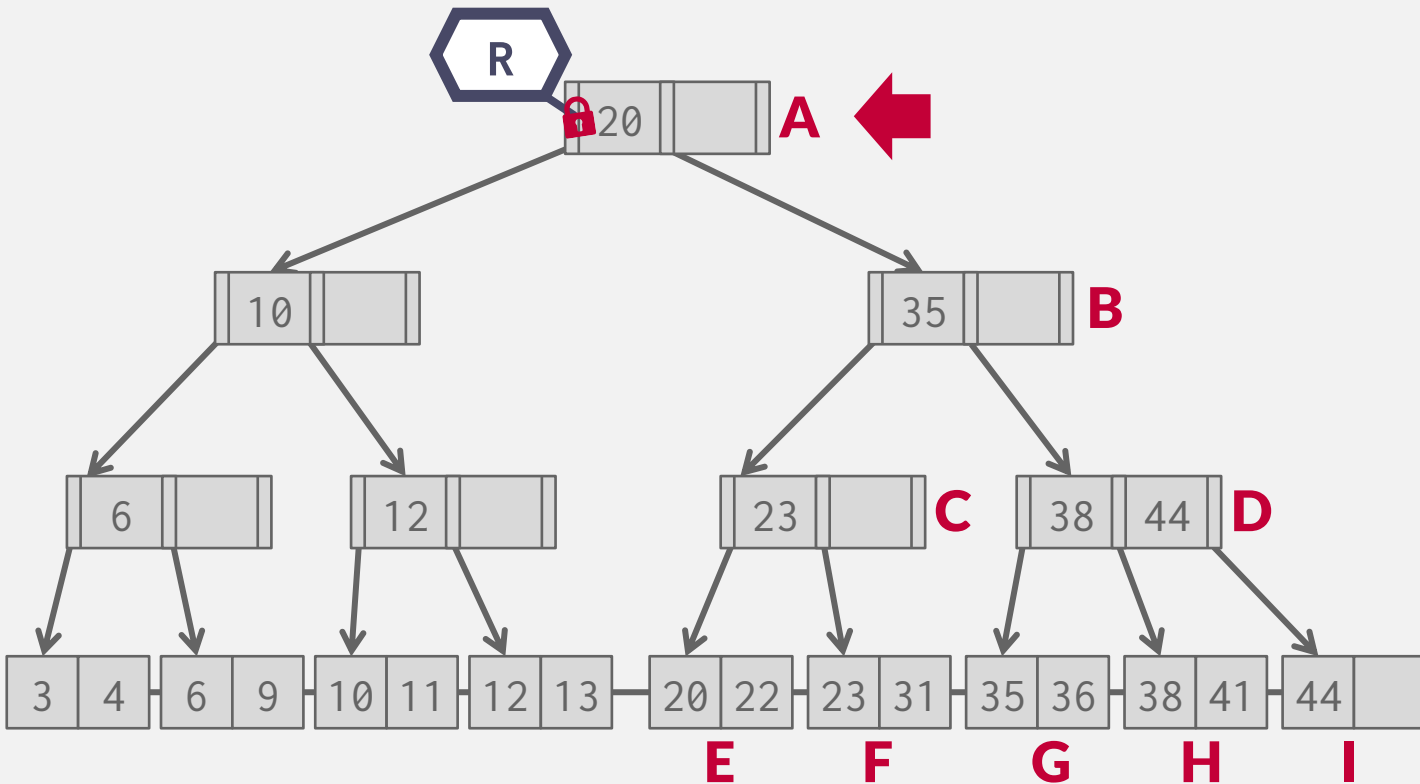
Find: Start at root and traverse down the tree:

- Acquire **R** latch on child,
- Then unlatch parent.
- Repeat until we reach the leaf node.

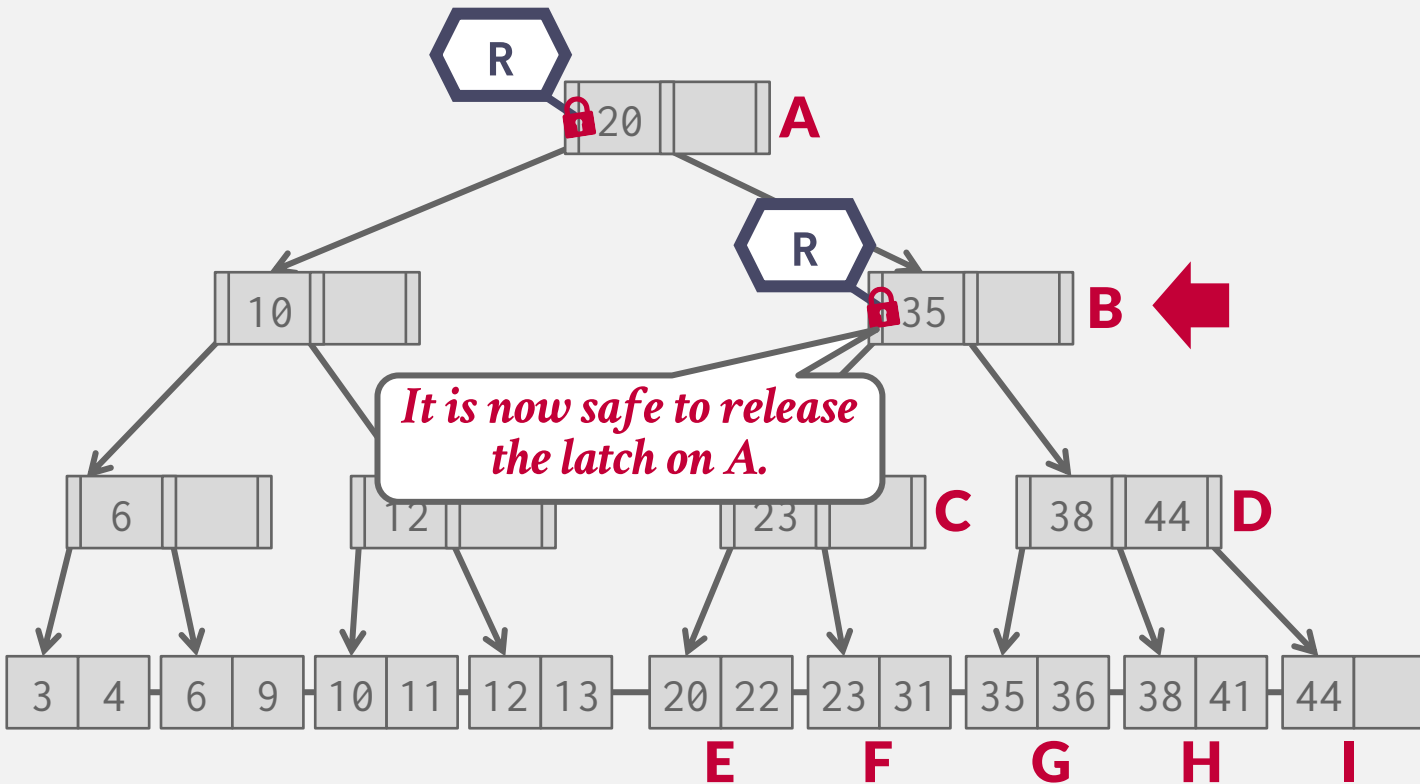
Insert/Delete: Start at root and go down, obtaining **W** latches as needed. Once child is latched, check if it is safe:

- If child is safe, release all latches on ancestors

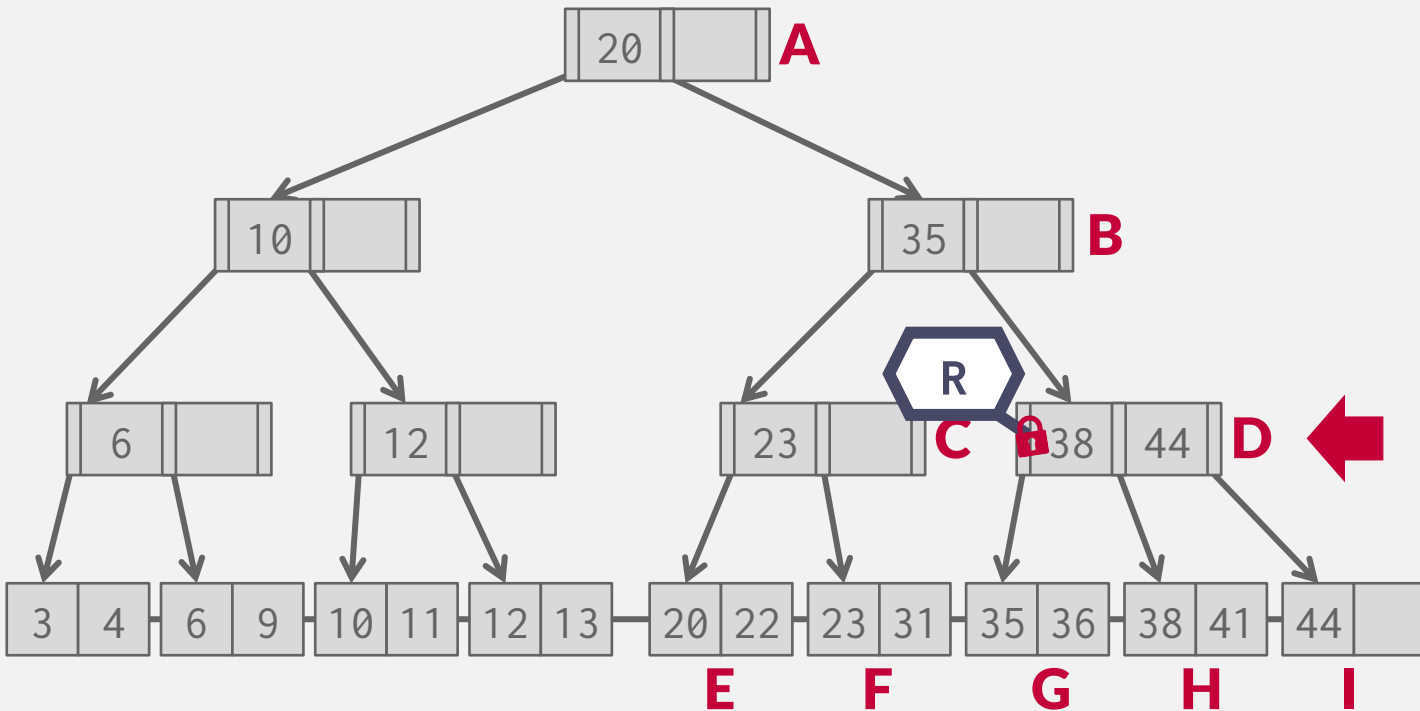
EXAMPLE #1 - FIND 38



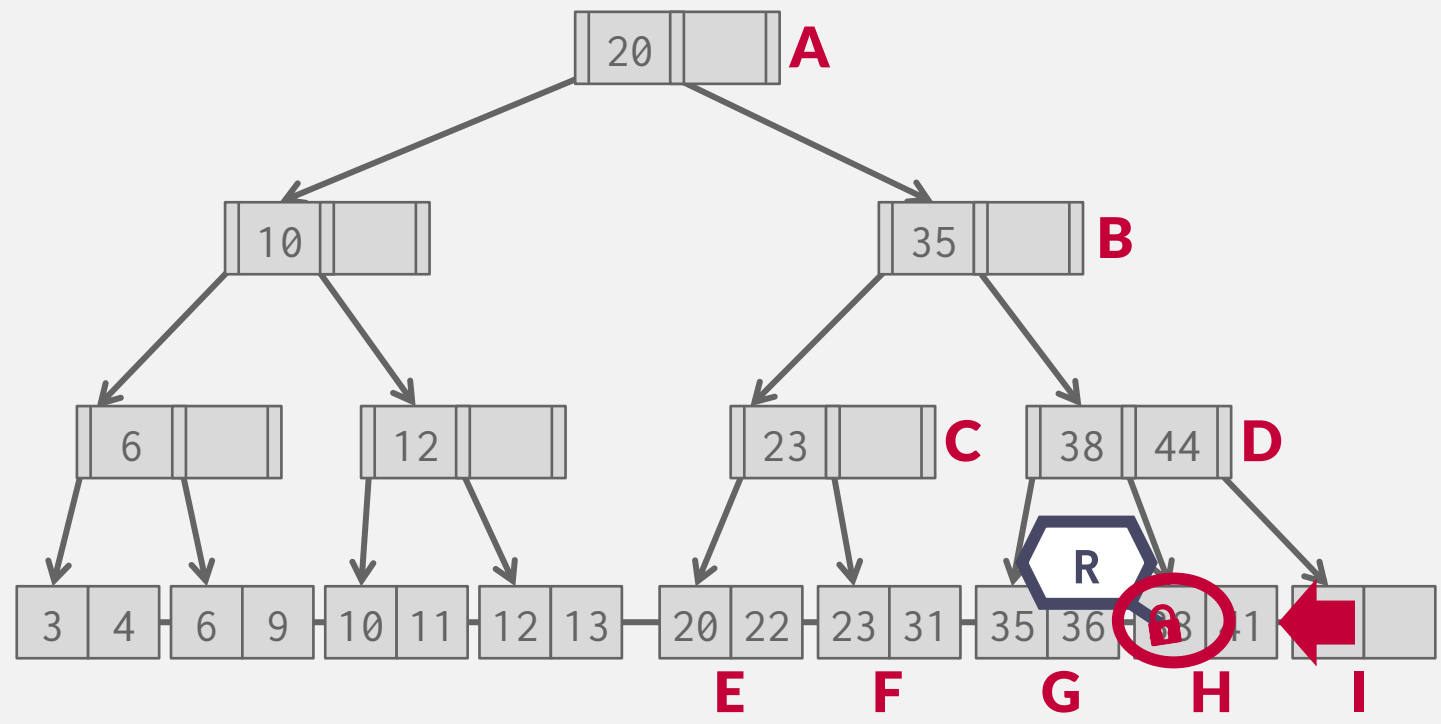
EXAMPLE #1 - FIND 38



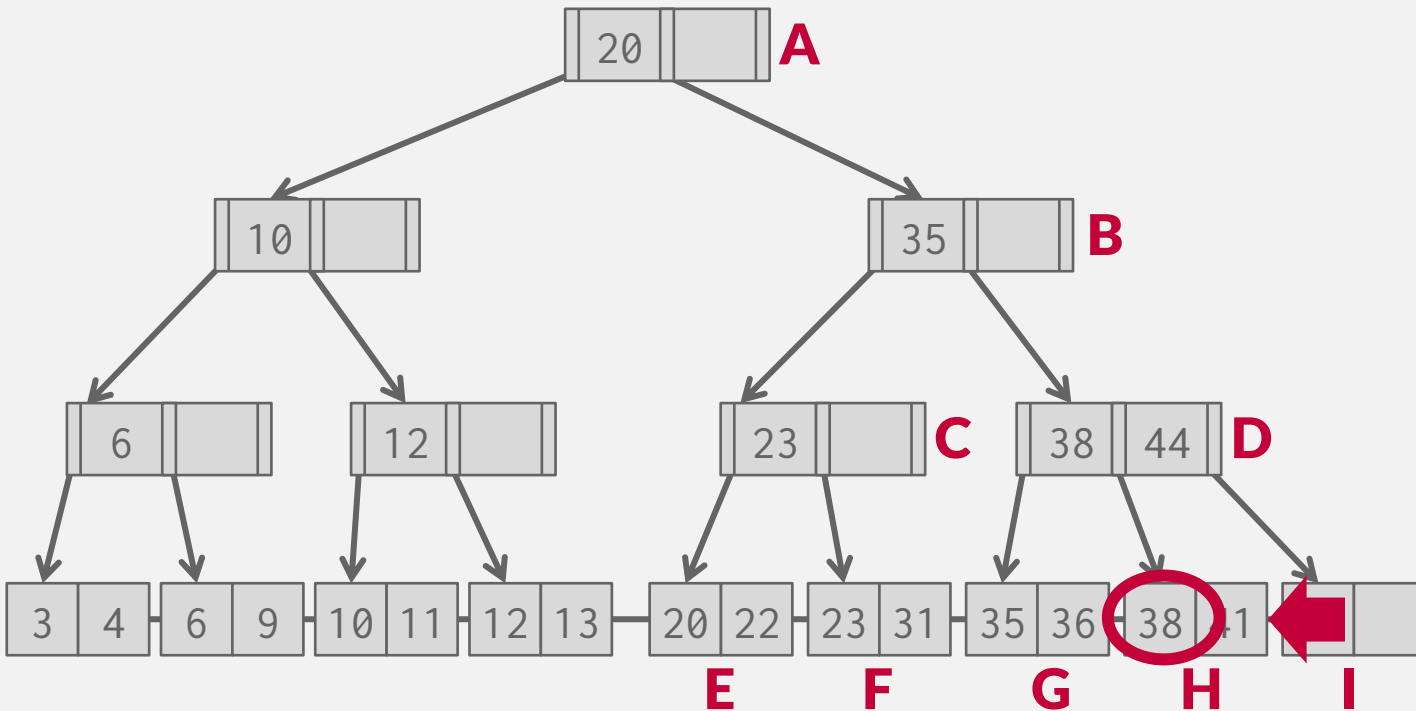
EXAMPLE #1 - FIND 38



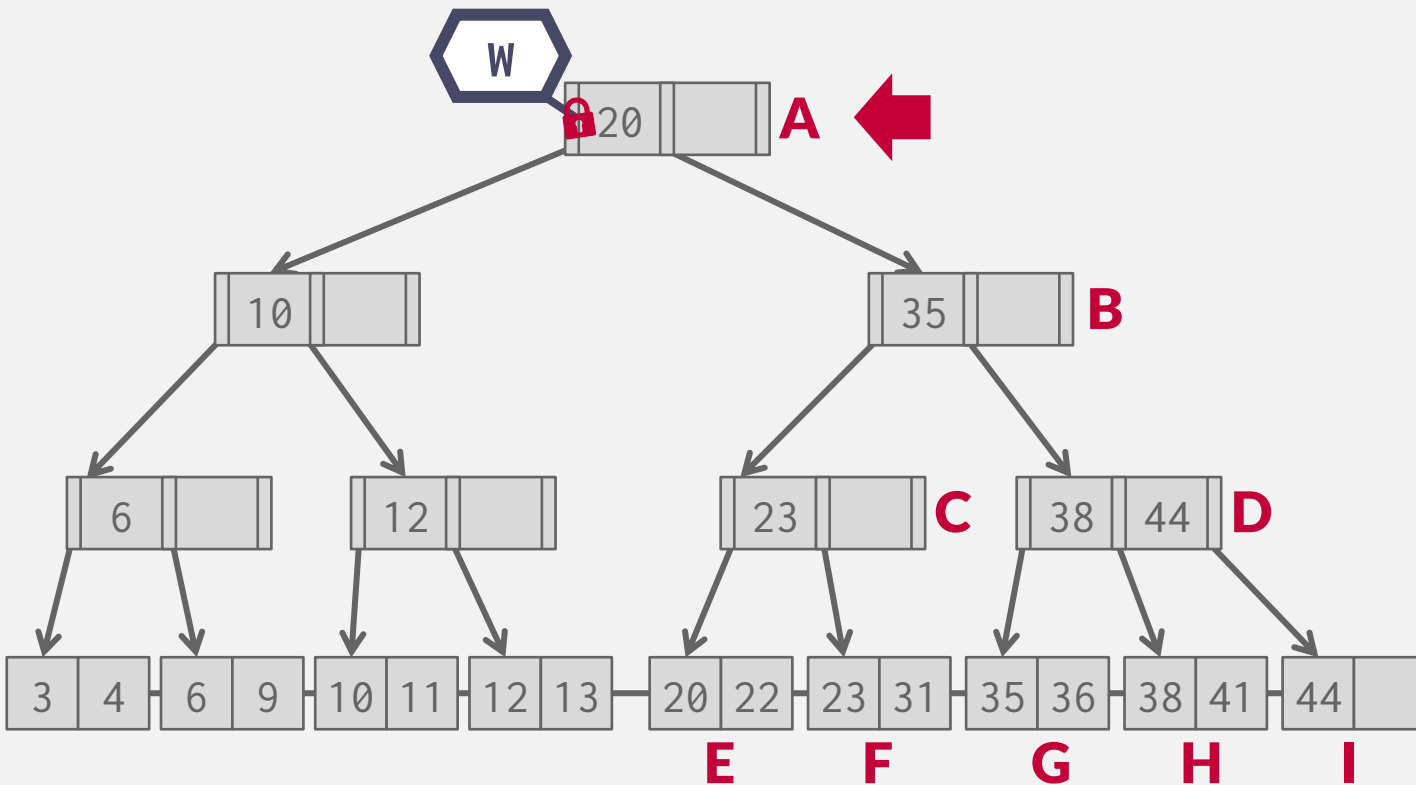
EXAMPLE #1 - FIND 38



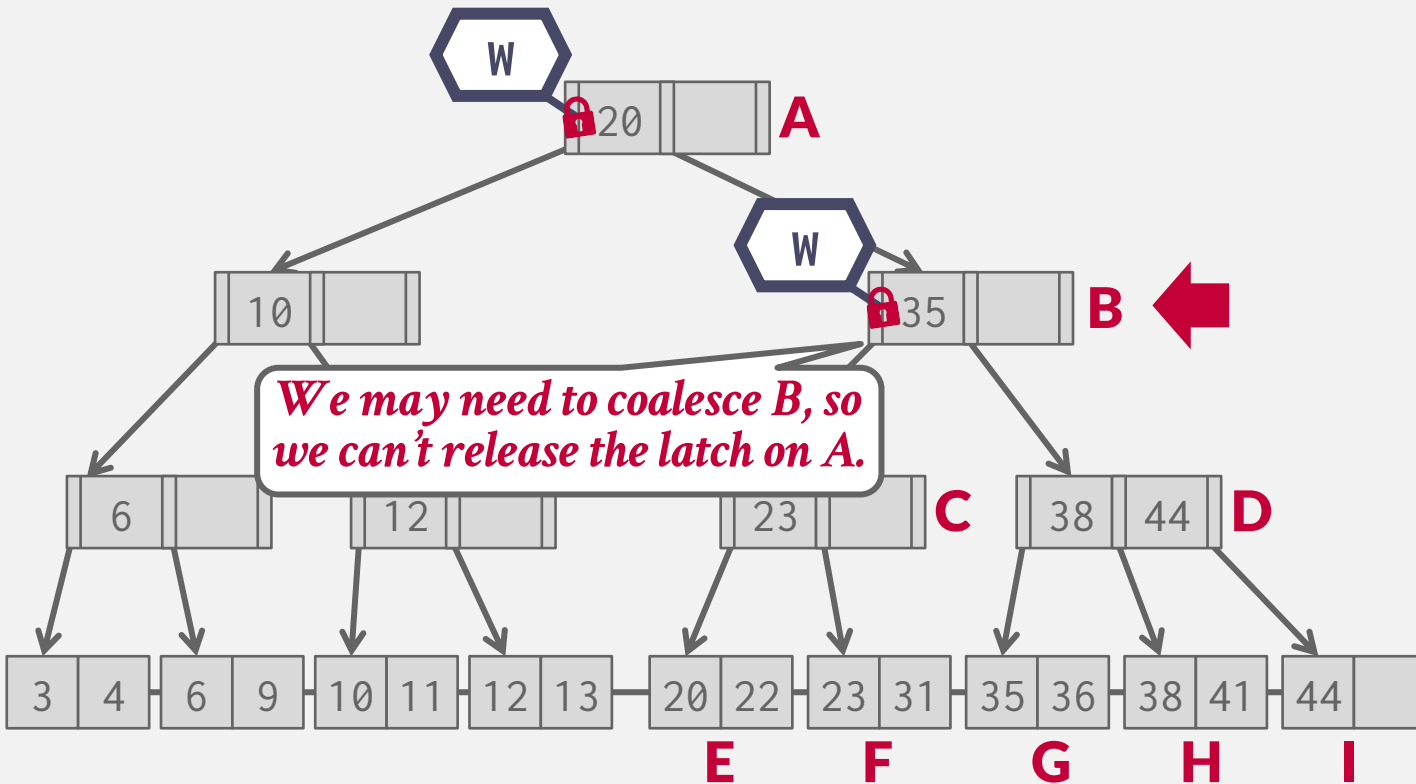
EXAMPLE #1 - FIND 38



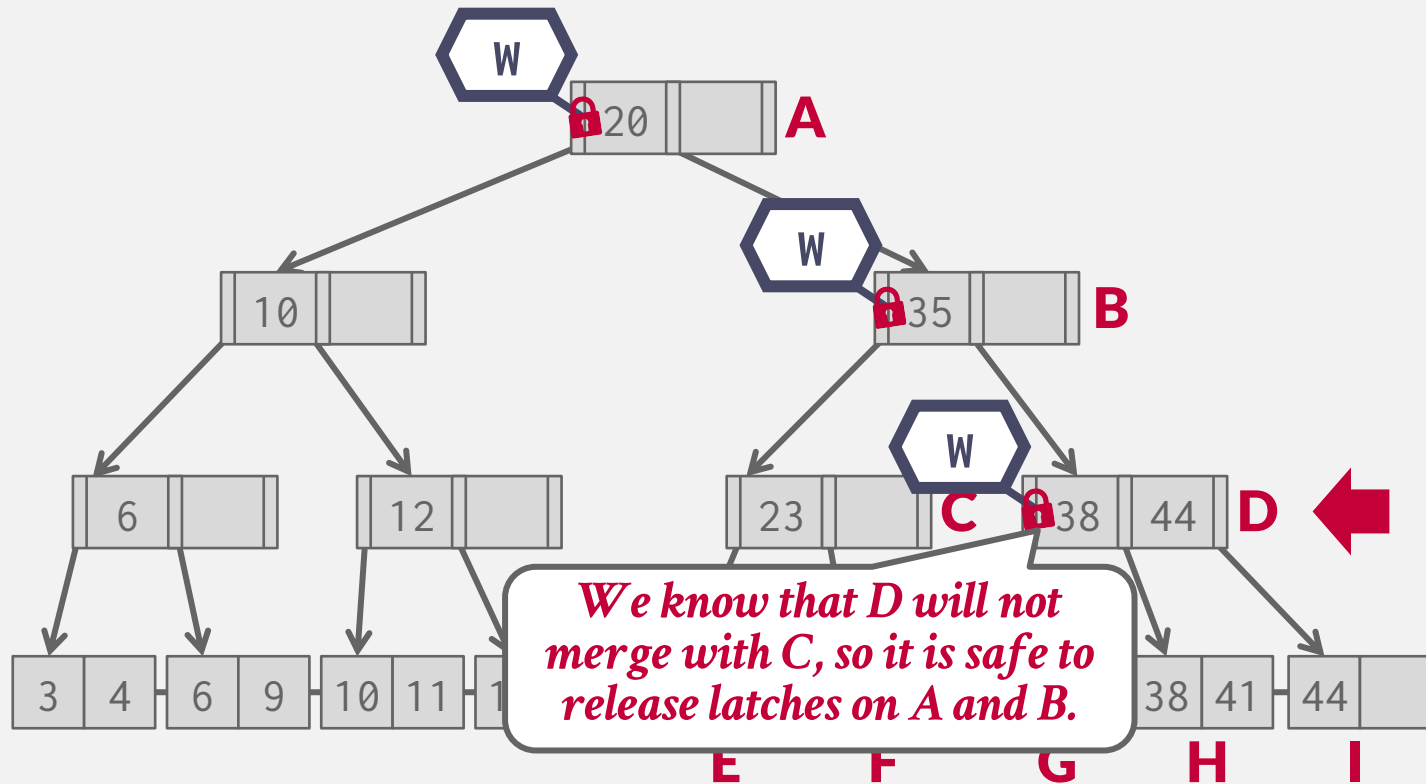
EXAMPLE #2 - DELETE 38



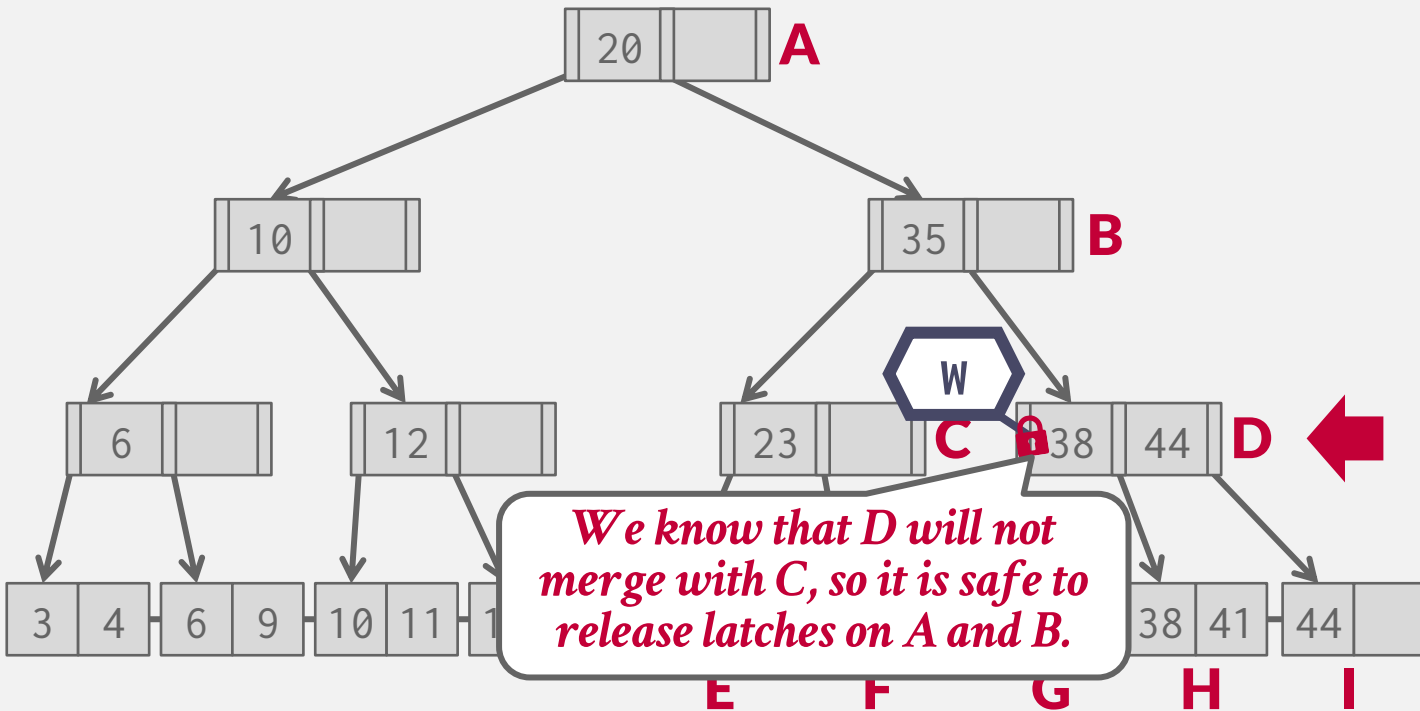
EXAMPLE #2 - DELETE 38



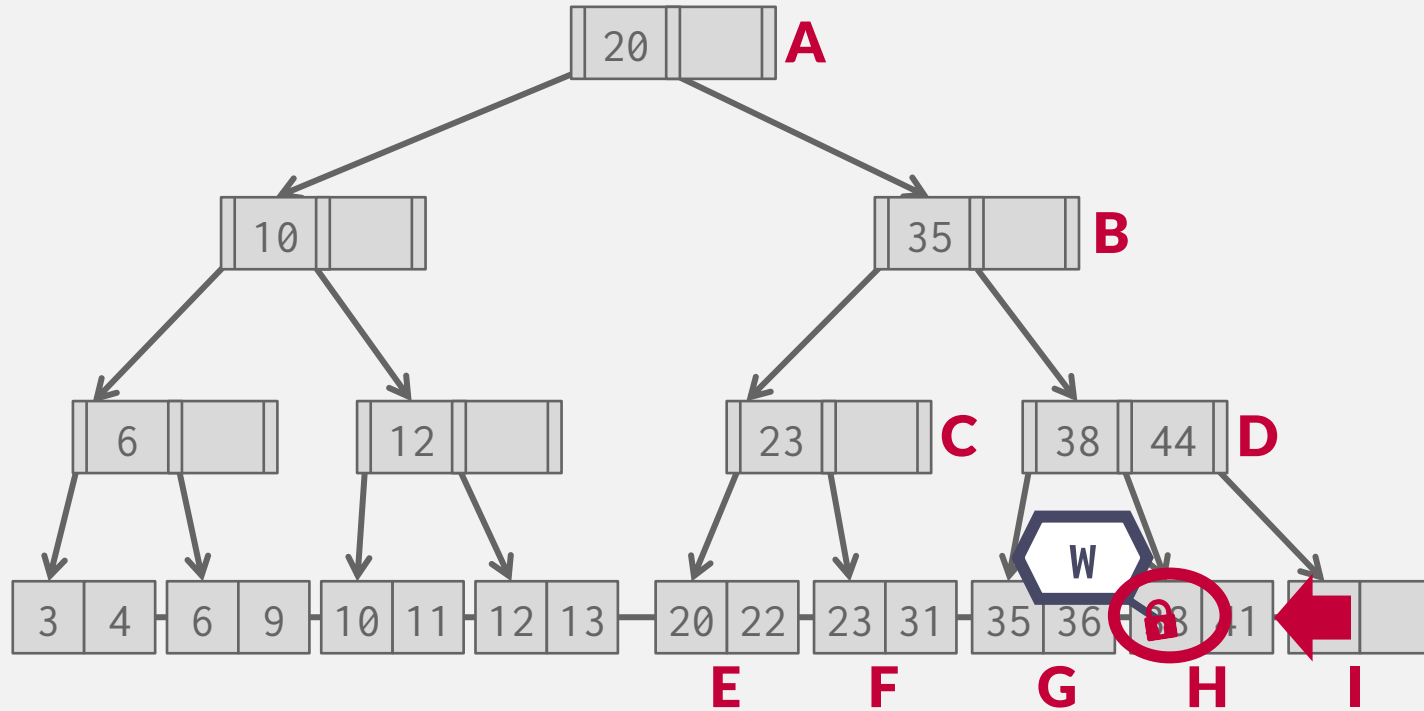
EXAMPLE #2 - DELETE 38



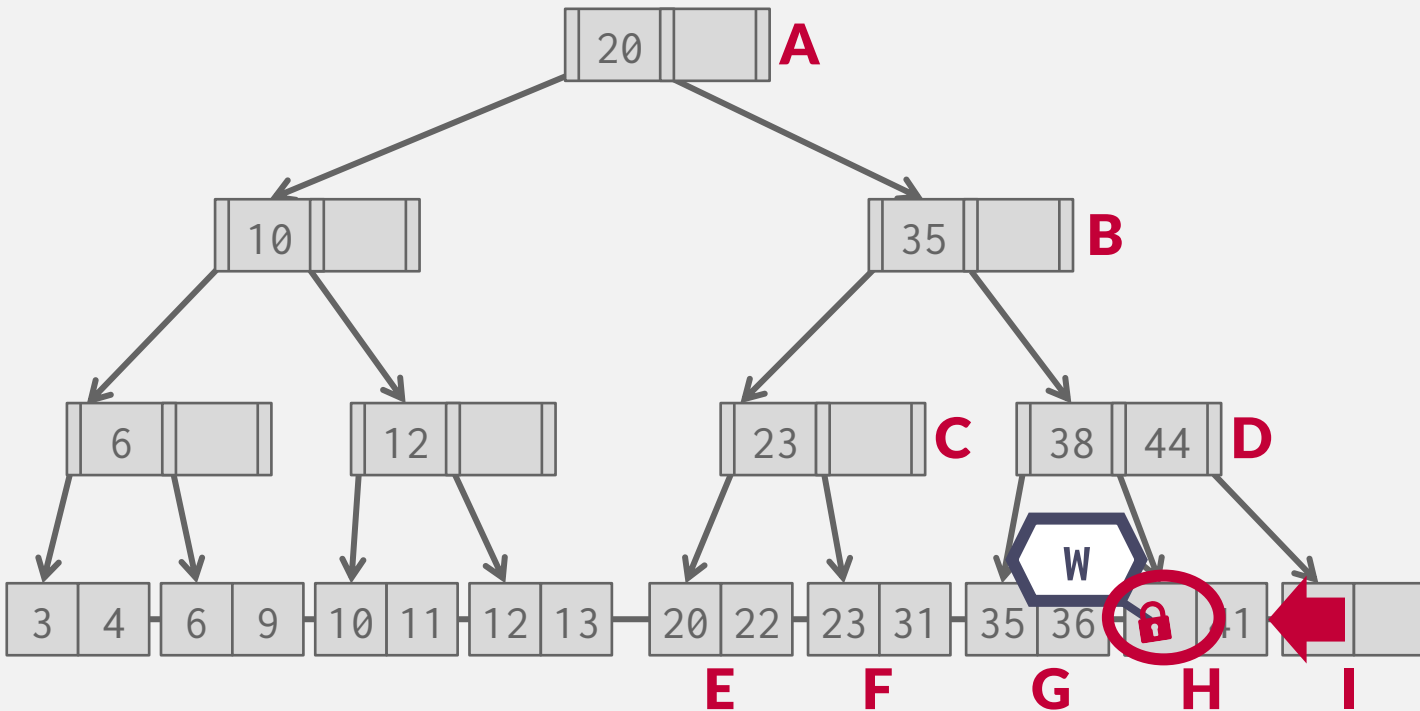
EXAMPLE #2 - DELETE 38



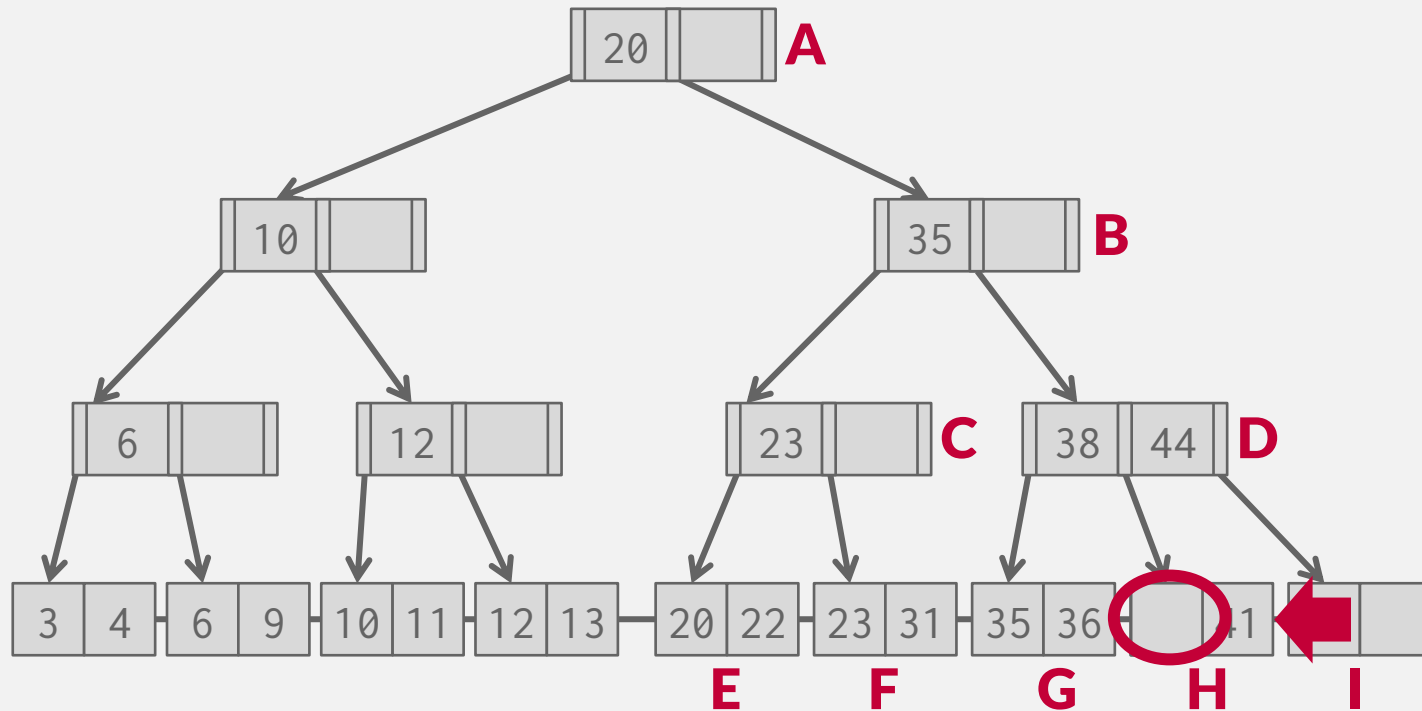
EXAMPLE #2 - DELETE 38



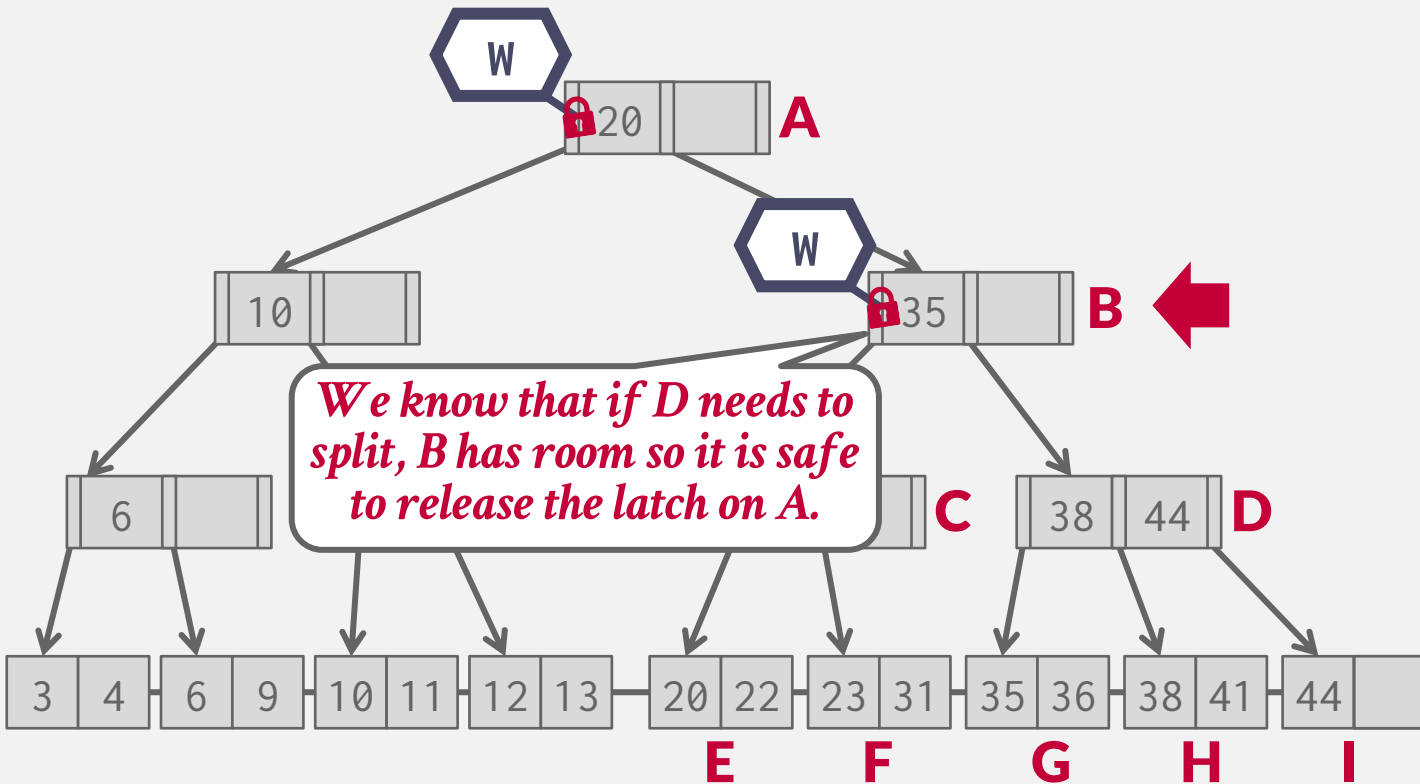
EXAMPLE #2 - DELETE 38



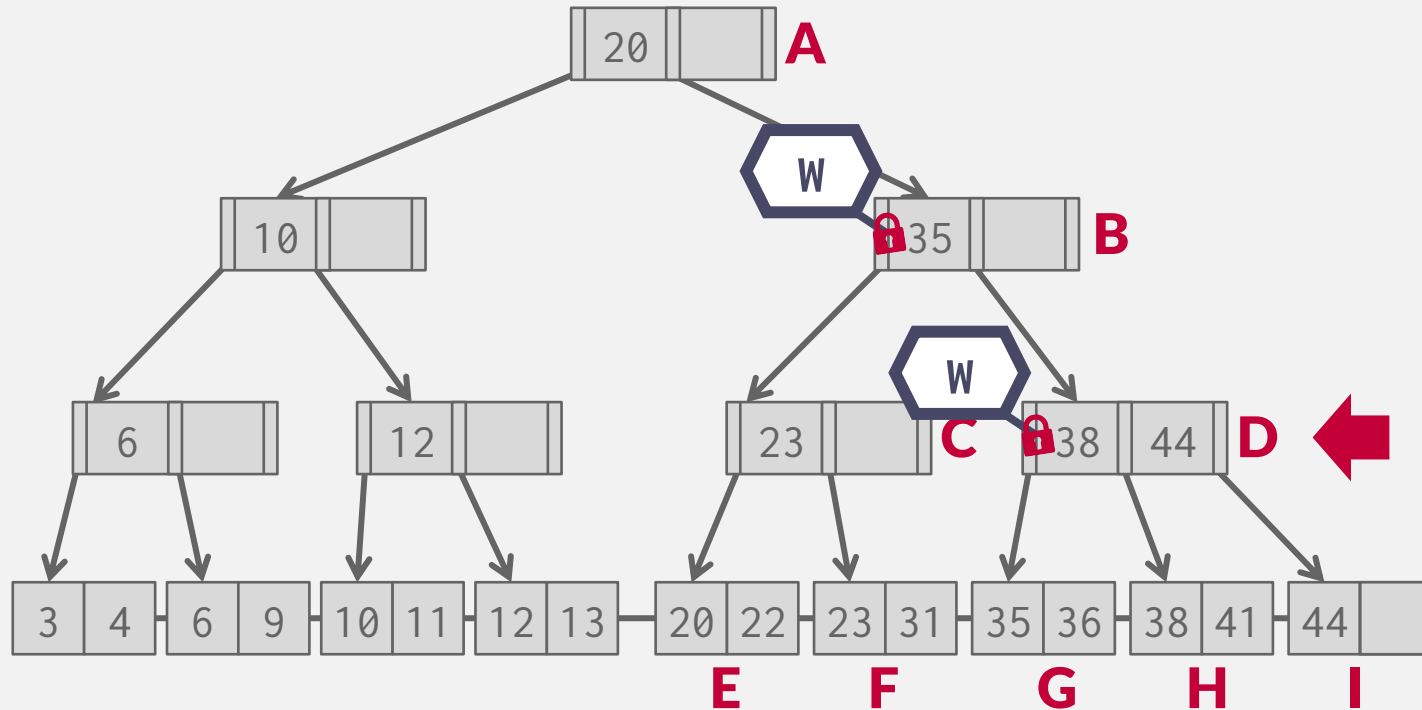
EXAMPLE #2 - DELETE 38



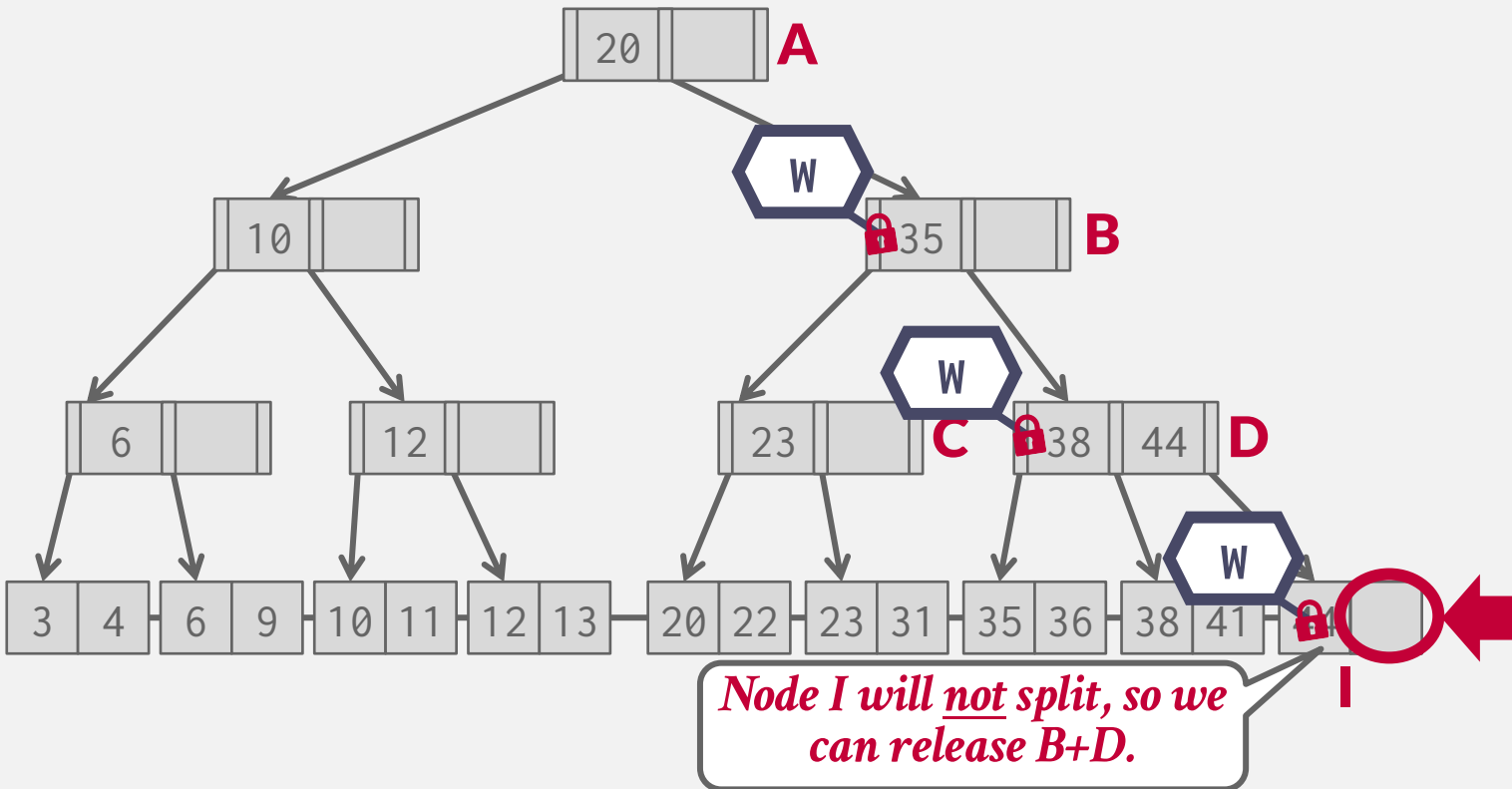
EXAMPLE #3 - INSERT 45



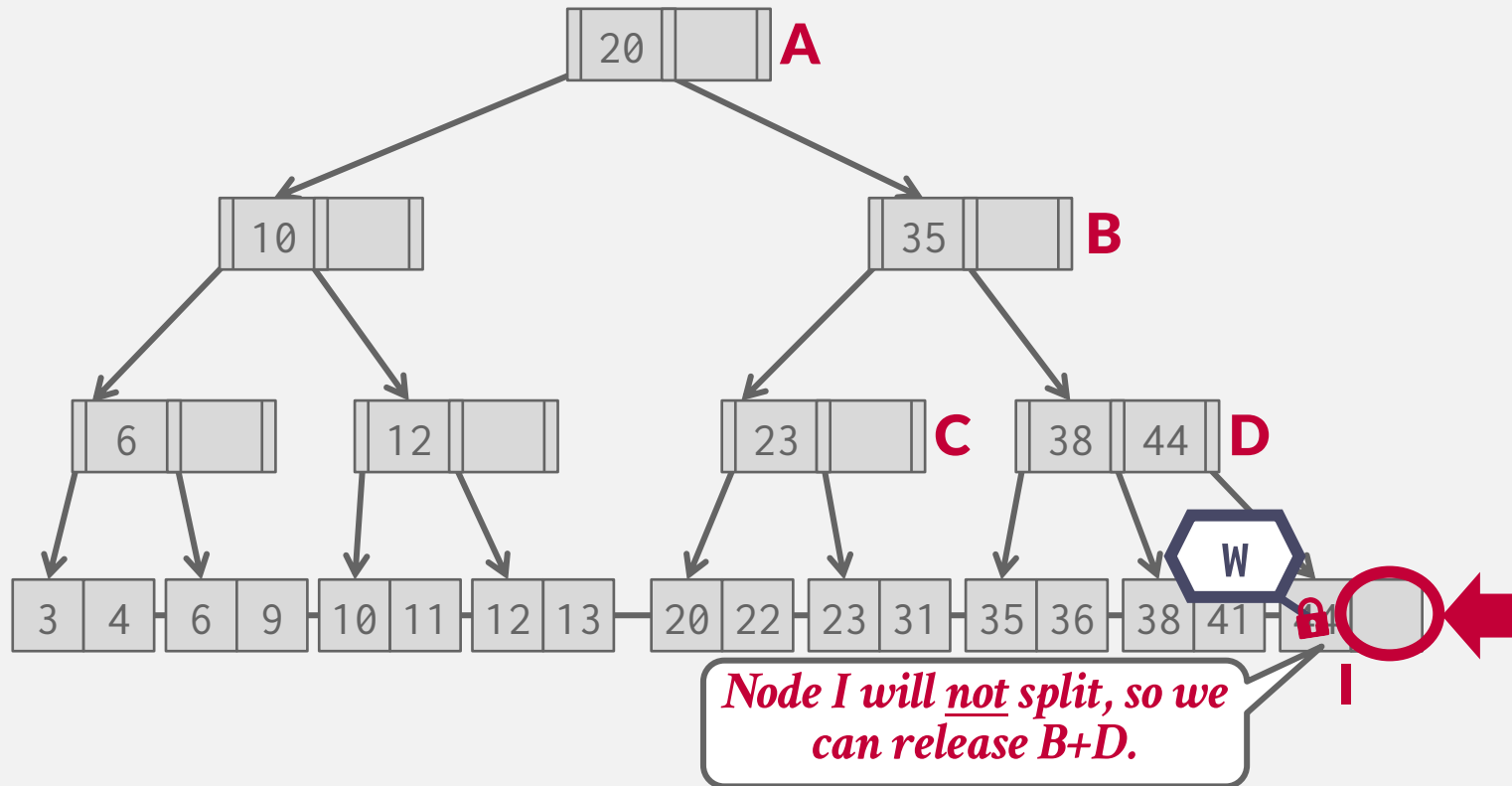
EXAMPLE #3 - INSERT 45



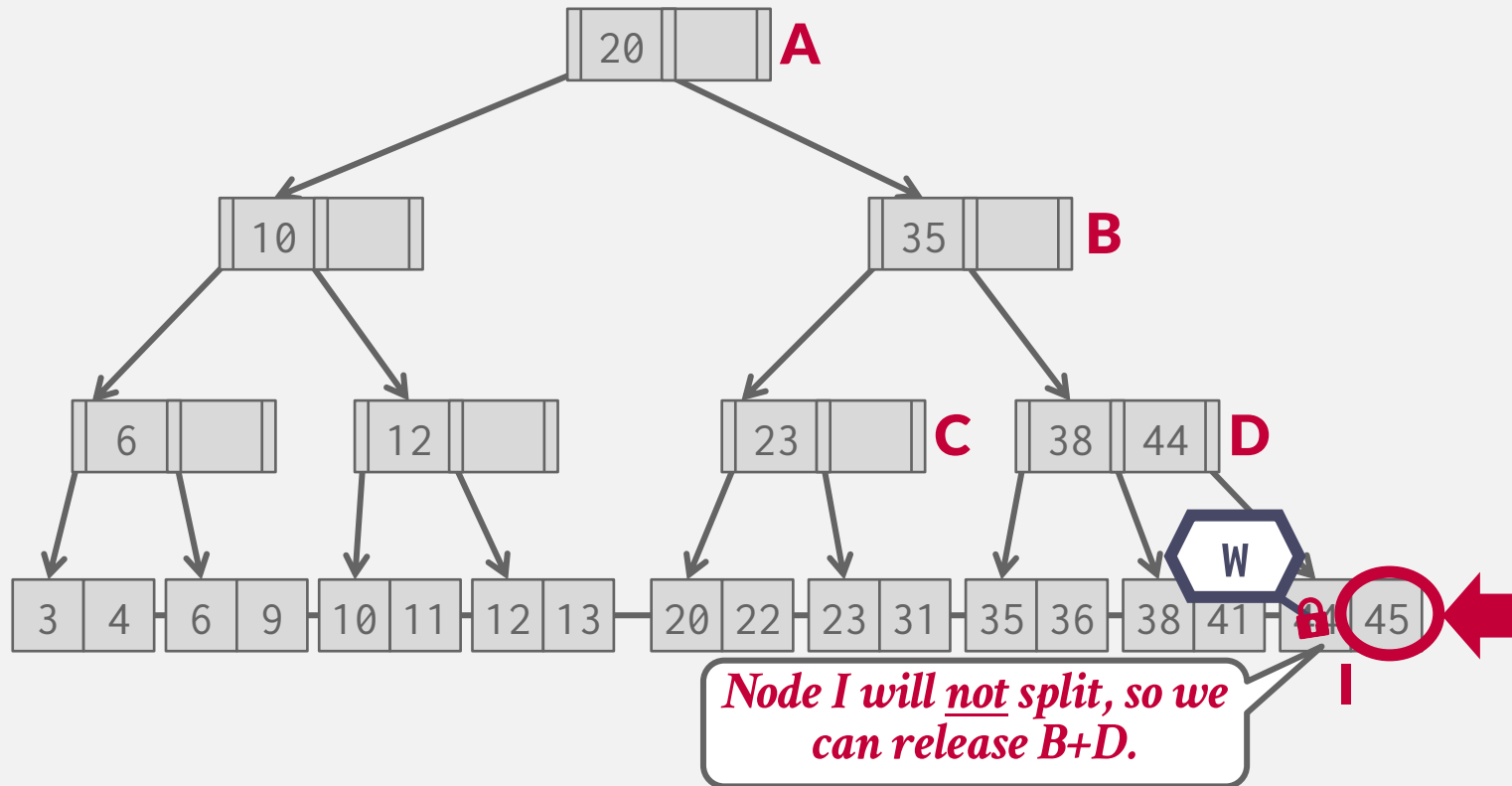
EXAMPLE #3 - INSERT 45



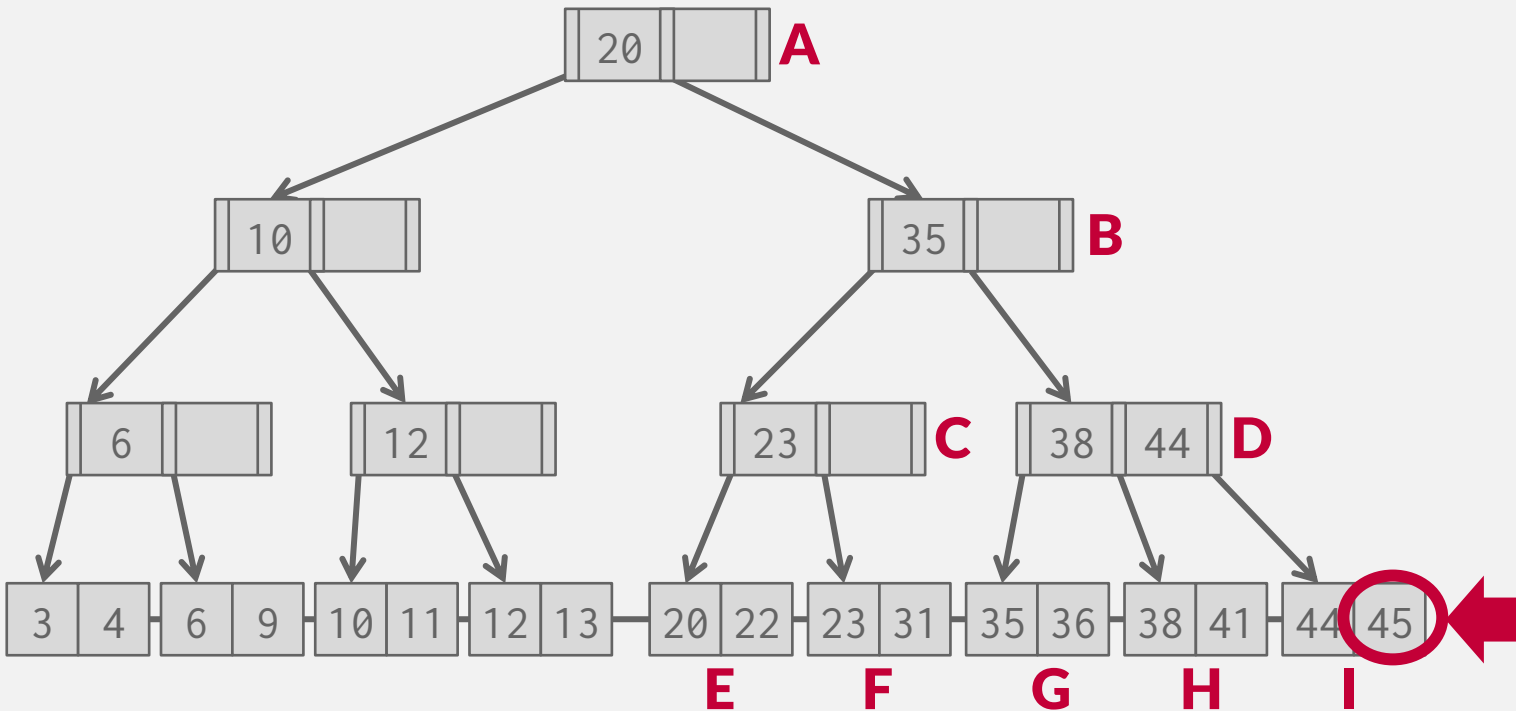
EXAMPLE #3 - INSERT 45



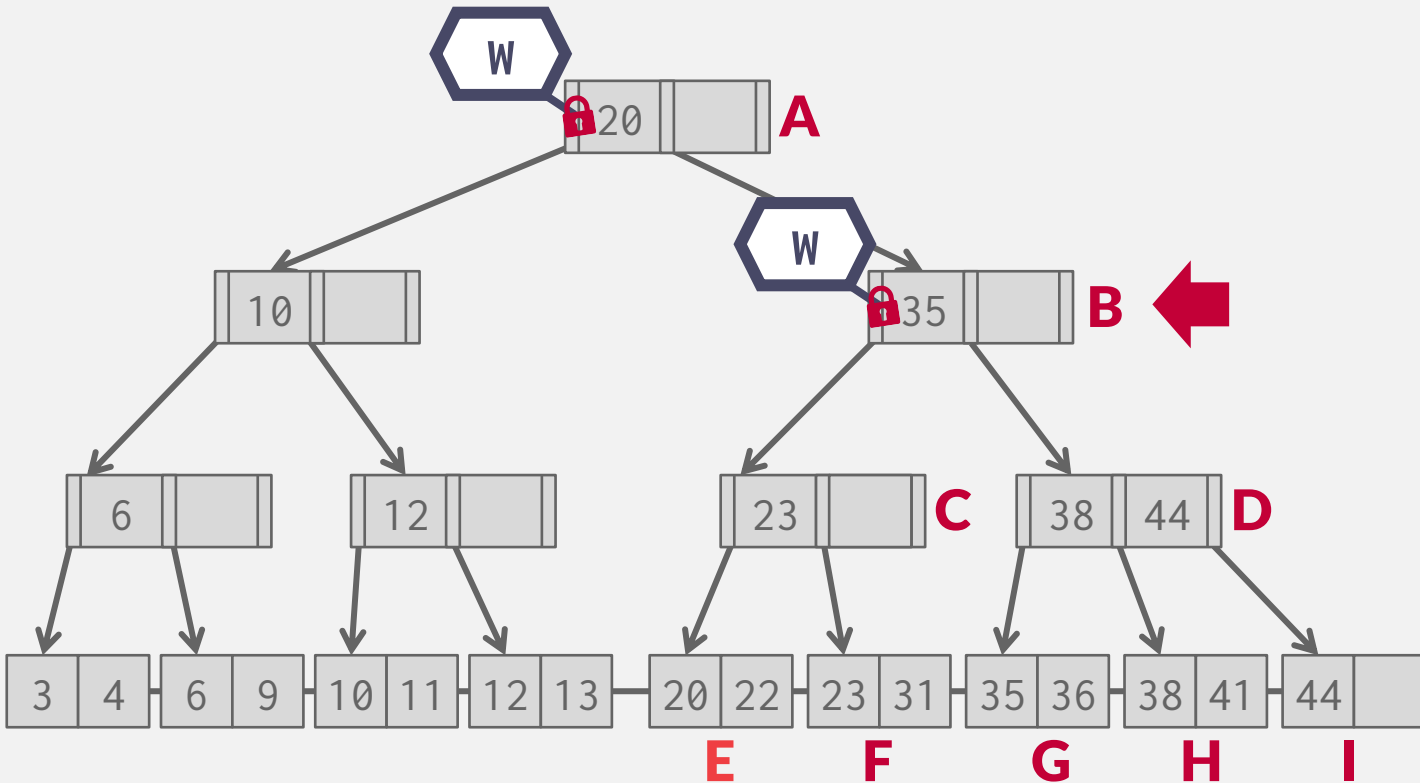
EXAMPLE #3 - INSERT 45



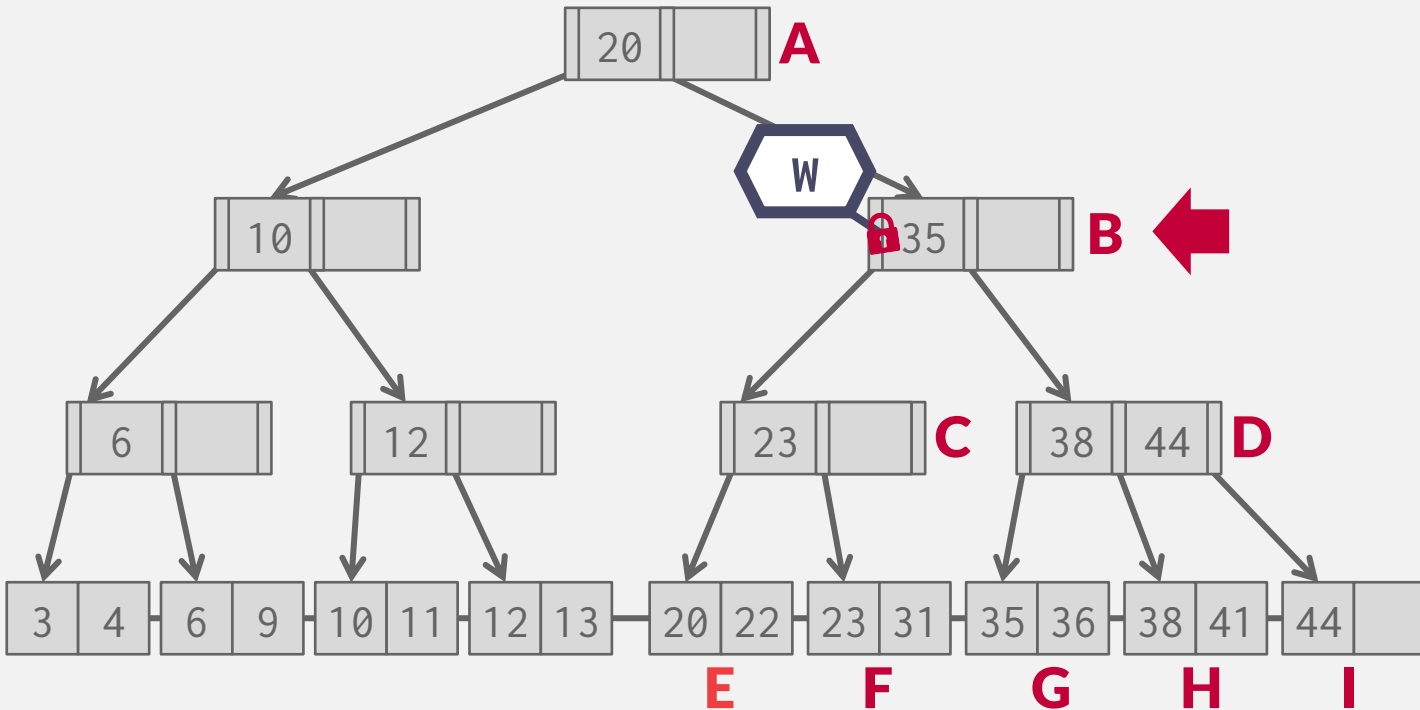
EXAMPLE #3 - INSERT 45



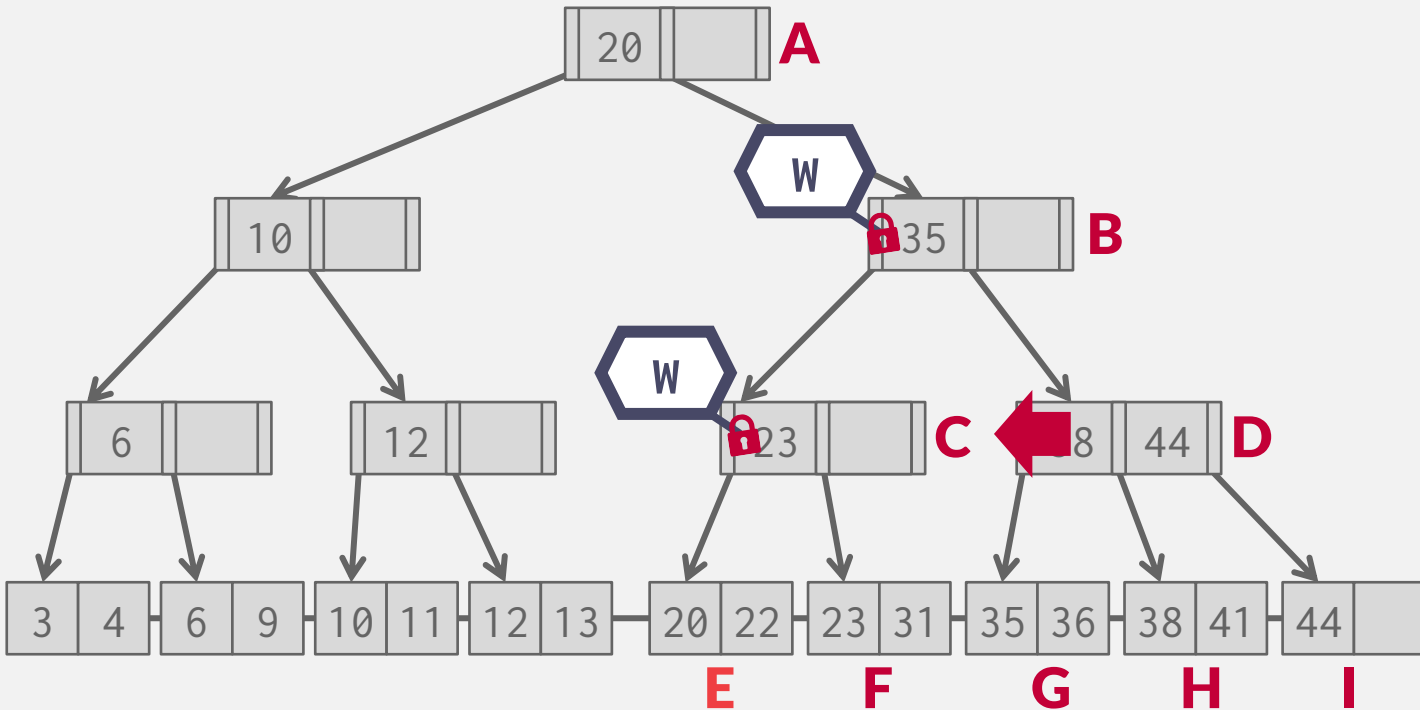
EXAMPLE #4 - INSERT 25



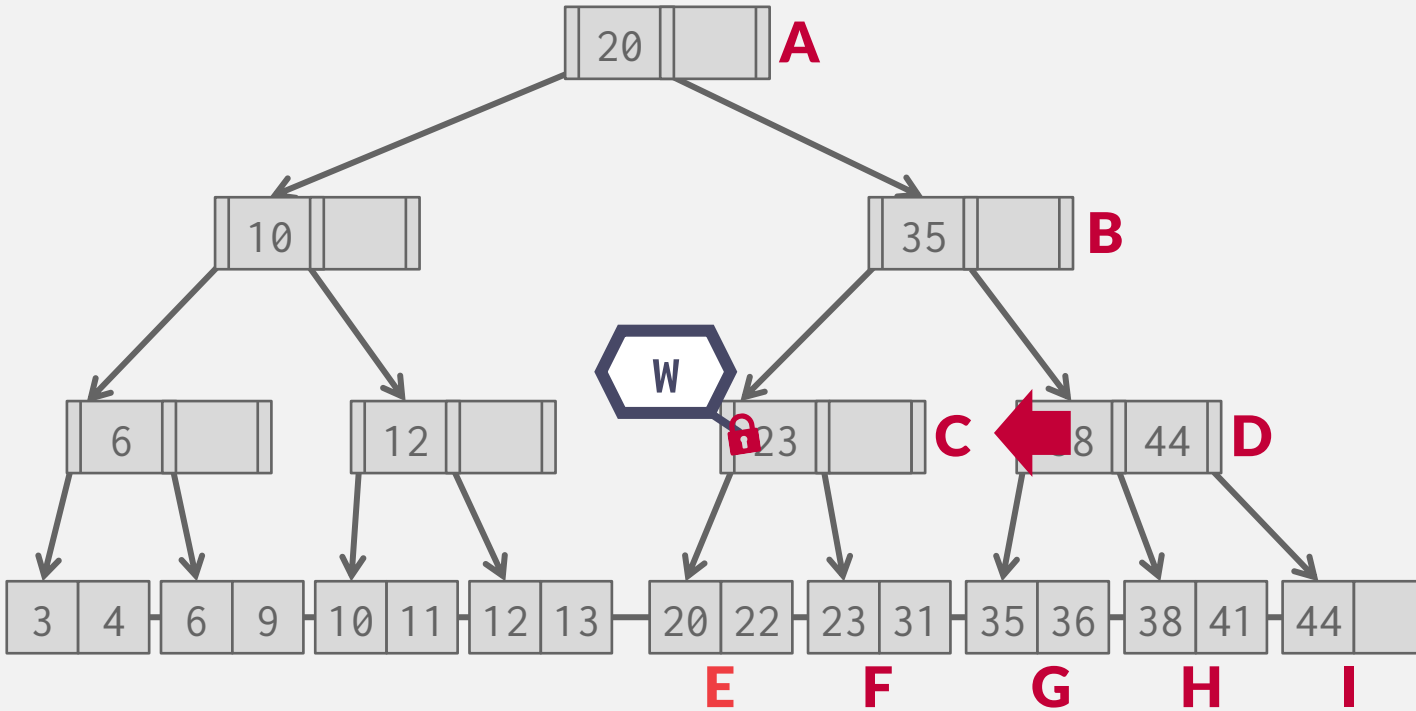
EXAMPLE #4 - INSERT 25



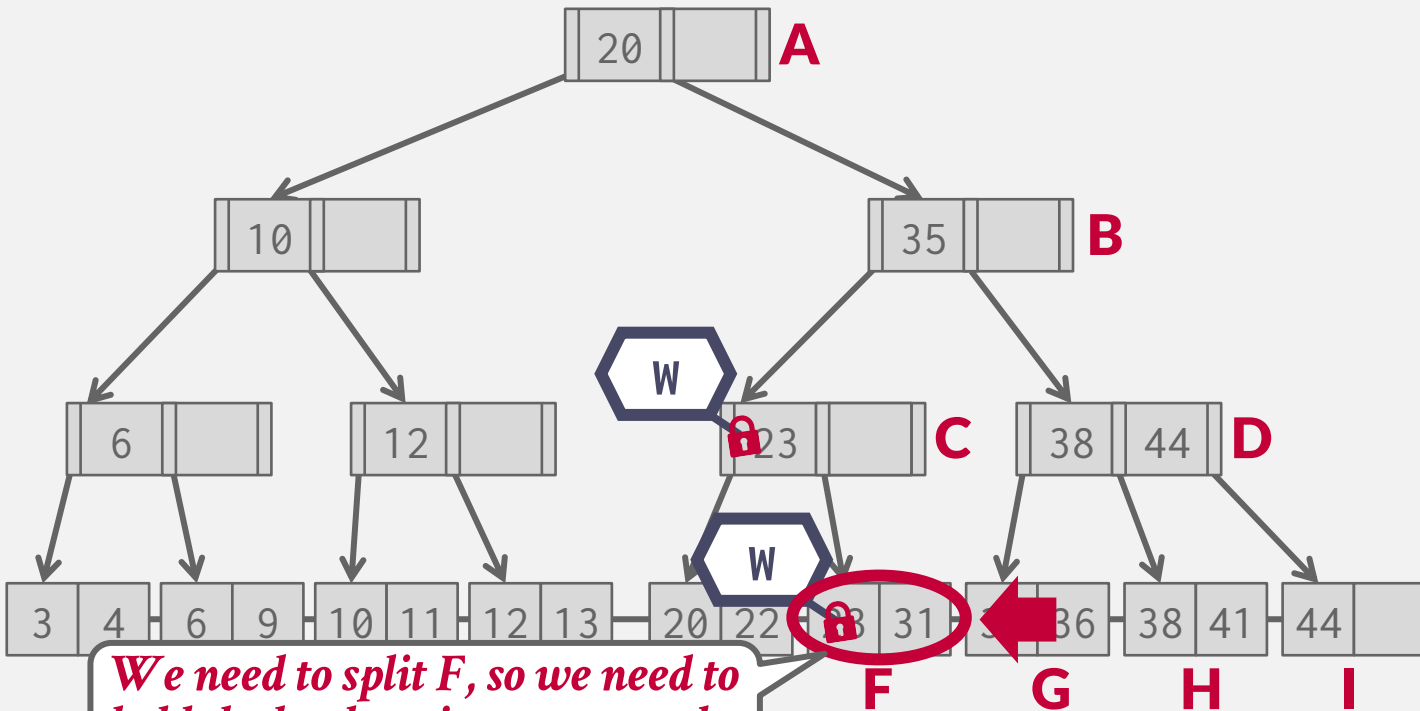
EXAMPLE #4 - INSERT 25



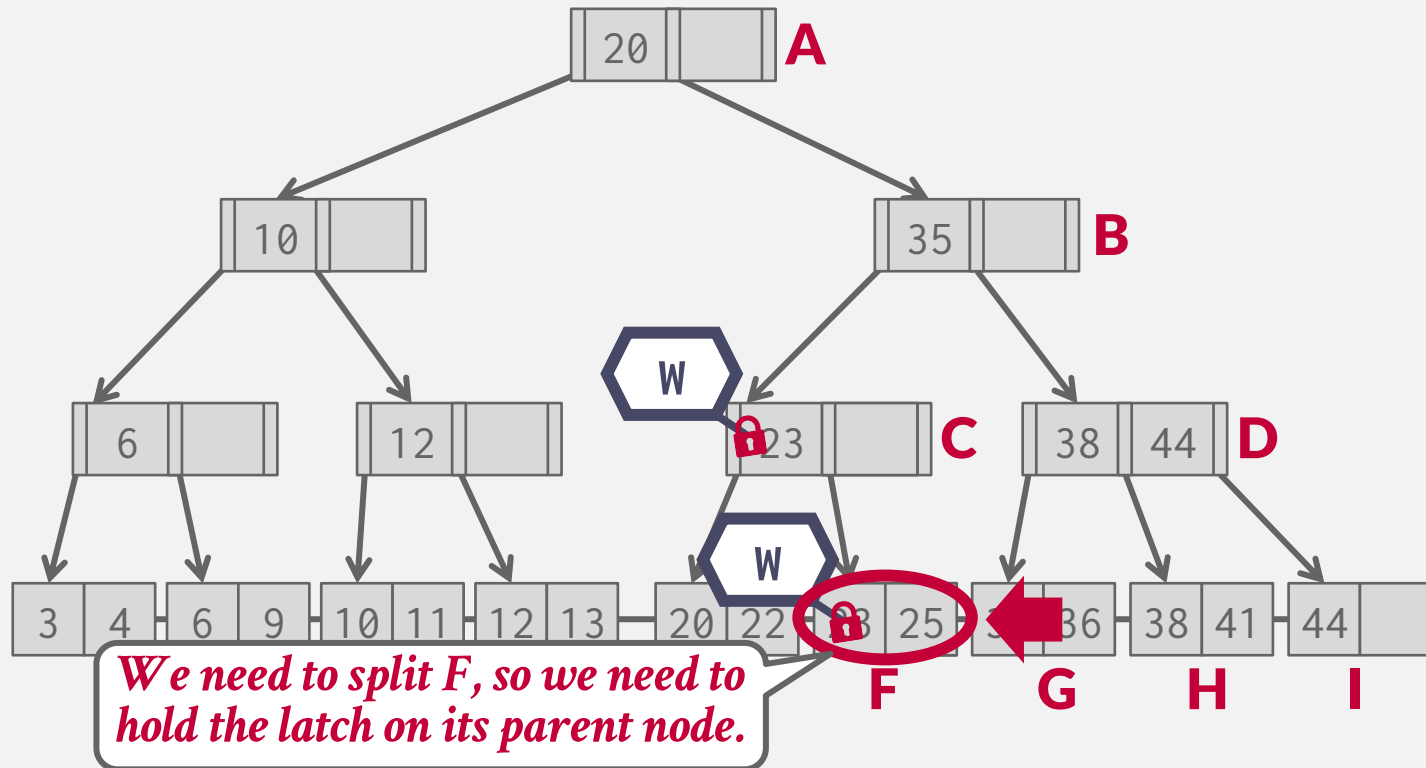
EXAMPLE #4 - INSERT 25



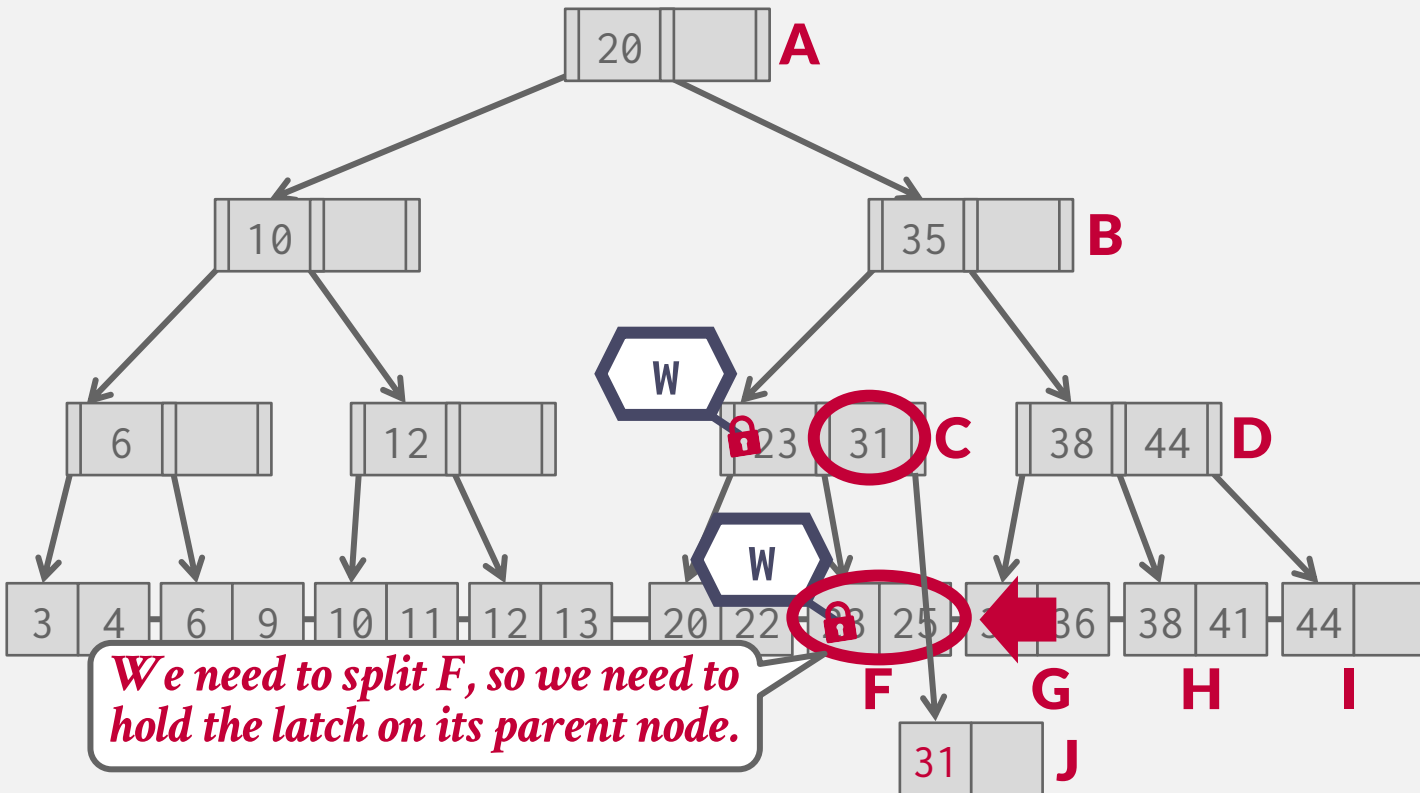
EXAMPLE #4 - INSERT 25



EXAMPLE #4 - INSERT 25

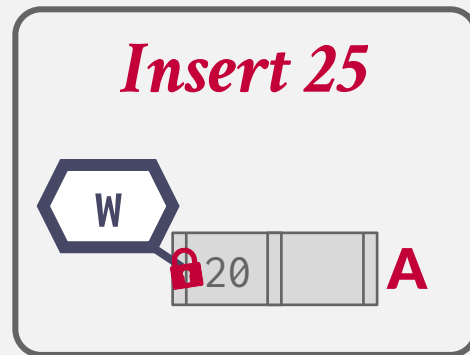
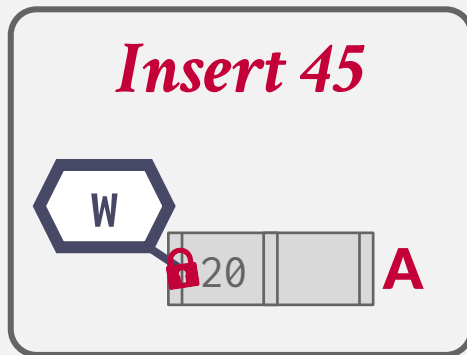
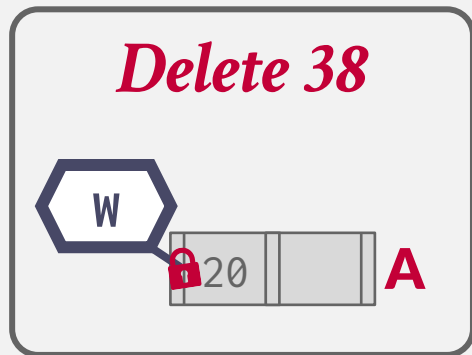


EXAMPLE #4 - INSERT 25



OBSERVATION

What was the first step that all the update examples did on the B+Tree?



Taking a write latch on the root every time becomes a bottleneck with higher concurrency.

BETTER LATCHING ALGORITHM

Most modifications to a B+Tree will not require a split or merge.

Instead of assuming that there will be a split/merge, optimistically traverse the tree using read latches.

If you guess wrong, repeat traversal with the pessimistic algorithm.

Acta Informatica 9, 1–21 (1977)


© by Springer-Verlag 1977

Concurrency of Operations on B-Trees

R. Bayer* and M. Schkolnick
IBM Research Laboratory, San José, CA 95193, USA

Summary. Concurrent operations on B-trees pose the problem of insuring that each operation can be carried out without interfering with other operations being performed simultaneously by other users. This problem can become critical if these structures are being used to support access paths, like indexes, to data base systems. In this case, serializing access to one of these indexes can create an unacceptable bottleneck for the entire system. Thus, there is a need for locking protocols that can assure integrity for each access while at the same time providing a maximum possible degree of concurrency. Another feature required from these protocols is that they be deadlock free, since the cost to resolve a deadlock may be high.

Recently, there has been some questioning on whether B-tree structures can support concurrent operations. In this paper, we examine the problem of concurrent access to B-trees. We present a deadlock free solution which can be tuned to specific requirements. An analysis is presented which allows the selection of parameters so as to satisfy these requirements.

The solution presented here uses simple locking protocols. Thus, we conclude that B-trees can be used advantageously in a multi-user environment.

1. Introduction

In this paper, we examine the problem of concurrent access to indexes which are maintained as B-trees. This type of organization was introduced by Bayer and McCreight [2] and some variants of it appear in Knuth [10] and Wedekind [13]. Performance studies of it were restricted to the single user environment. Recently, these structures have been examined for possible use in a multi-user (concurrent) environment. Some initial studies have been made about the feasibility of their use in this type of situation [1, 6], and [11].

An accessing schema which achieves a high degree of concurrency in using the index will be presented. The schema allows dynamic tuning to adapt its performance to the profile of the current set of users. Another property of the

* Permanent address: Institut für Informatik der Technischen Universität München, Arcisstr. 21, D-8000 München 2, Germany (Fed. Rep.)

BETTER LATCHING ALGORITHM

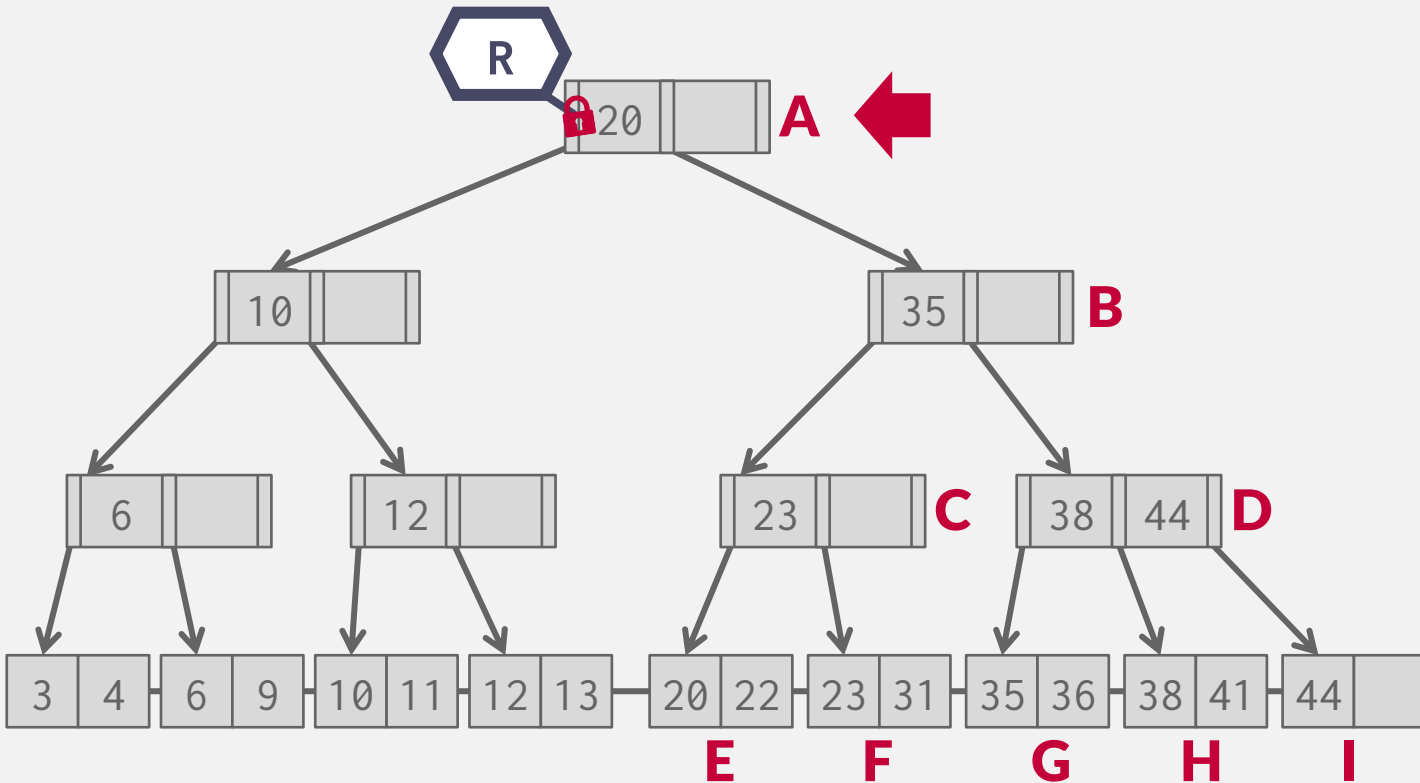
Search: Same as before.

Insert/Delete:

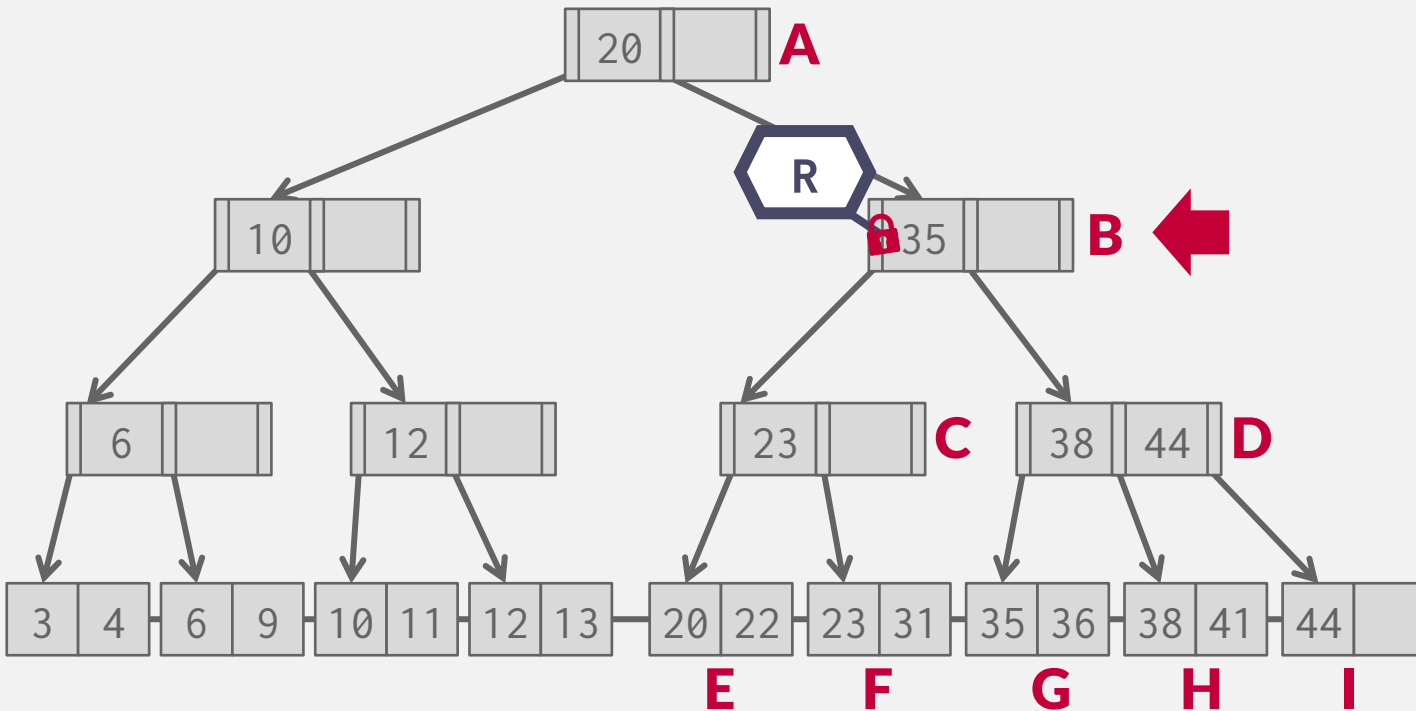
- Set latches as if for search, get to leaf, and set **W** latch on leaf.
- If leaf is not safe, release all latches, and restart thread using previous insert/delete protocol with write latches.

This approach optimistically assumes that only leaf node will be modified; if not, **R** latches set on the first pass to leaf are wasteful.

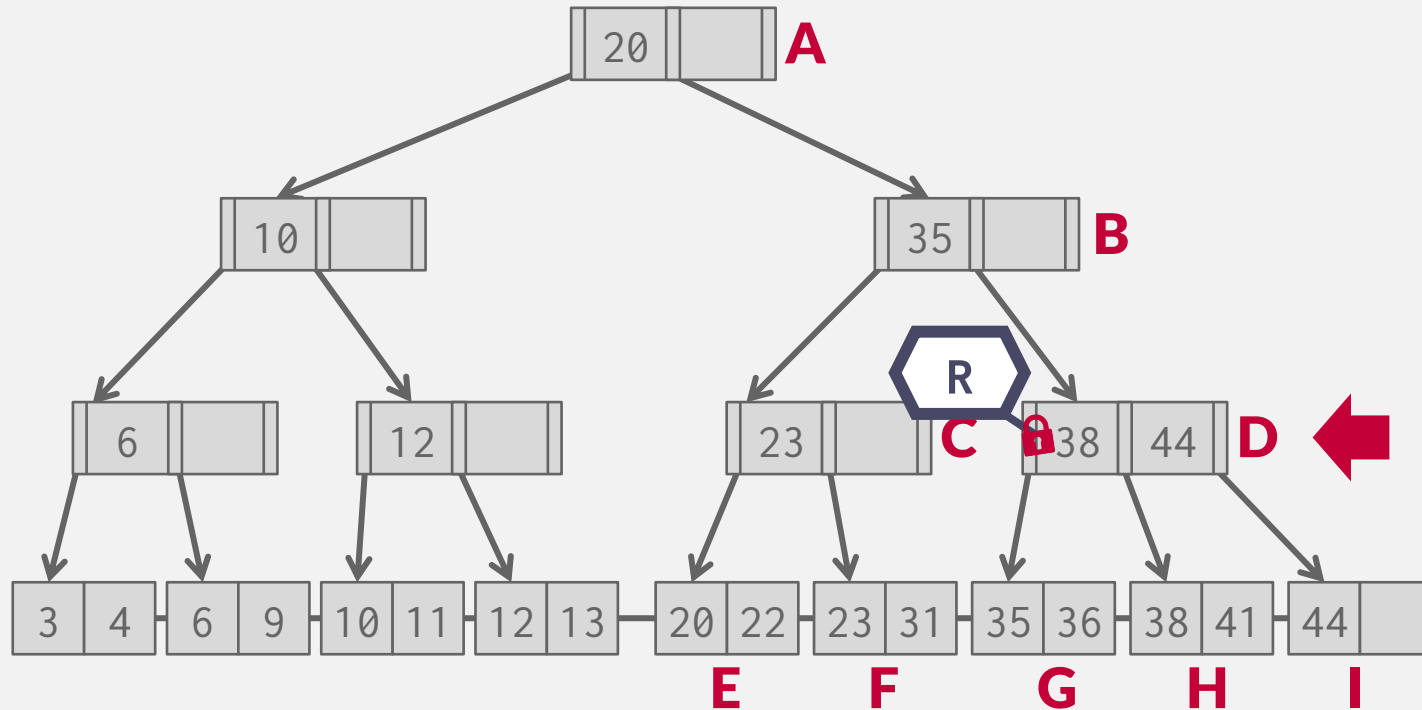
EXAMPLE #2 - DELETE 38



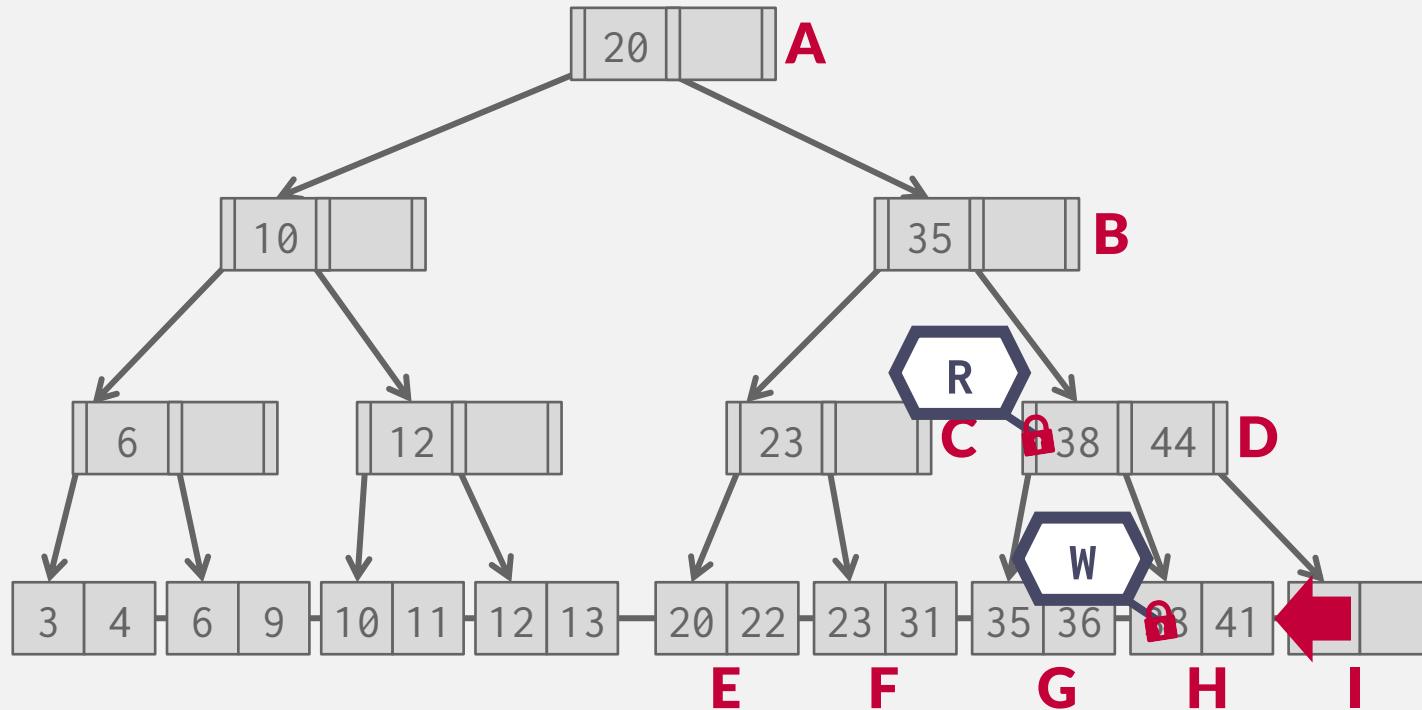
EXAMPLE #2 - DELETE 38



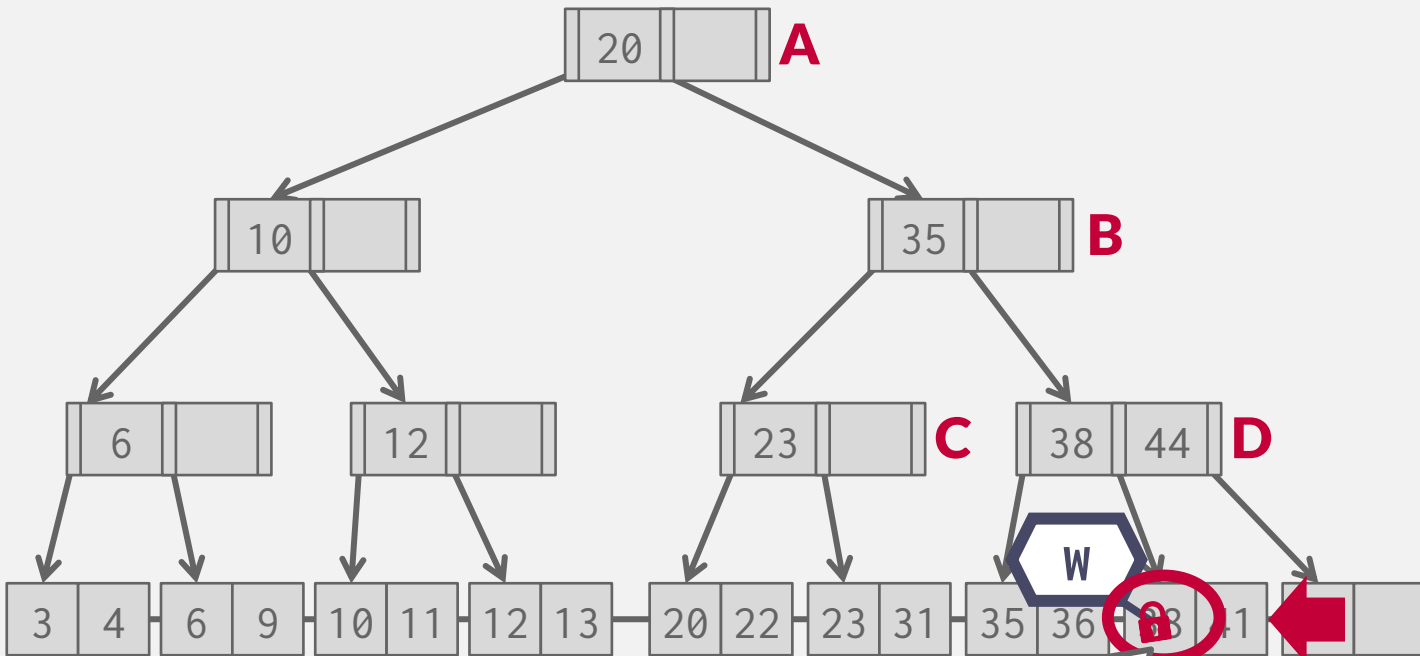
EXAMPLE #2 - DELETE 38



EXAMPLE #2 - DELETE 38

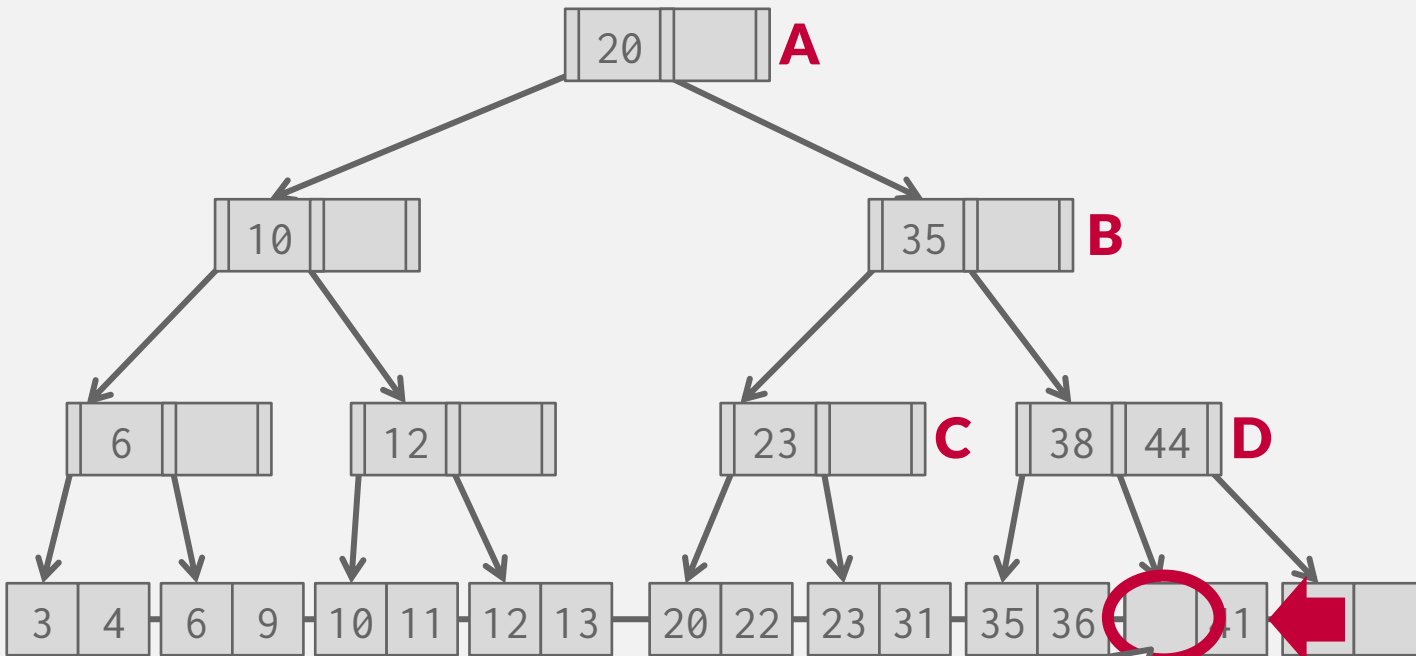


EXAMPLE #2 - DELETE 38

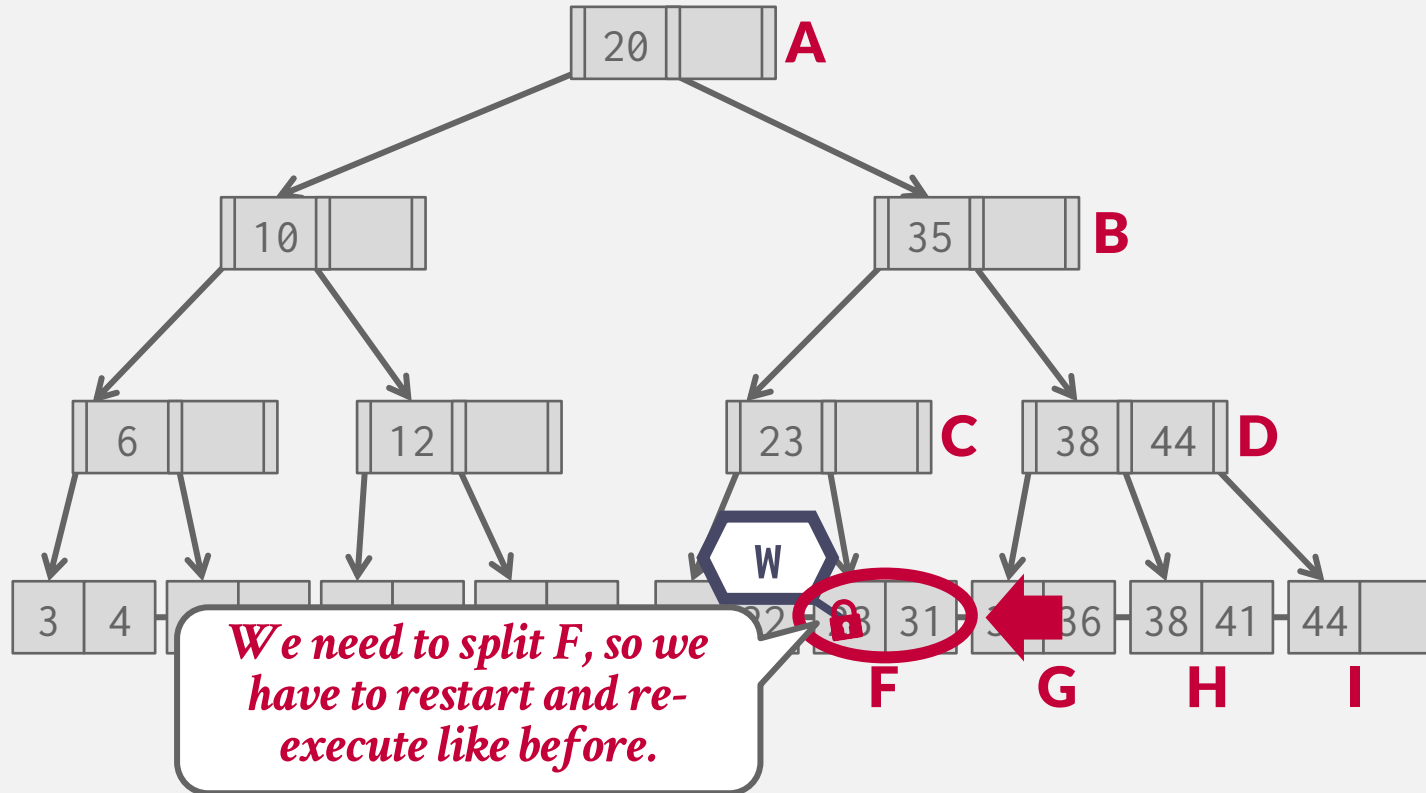


*Node H will not coalesce,
so we're safe!*

EXAMPLE #2 - DELETE 38



EXAMPLE #4 - INSERT 25



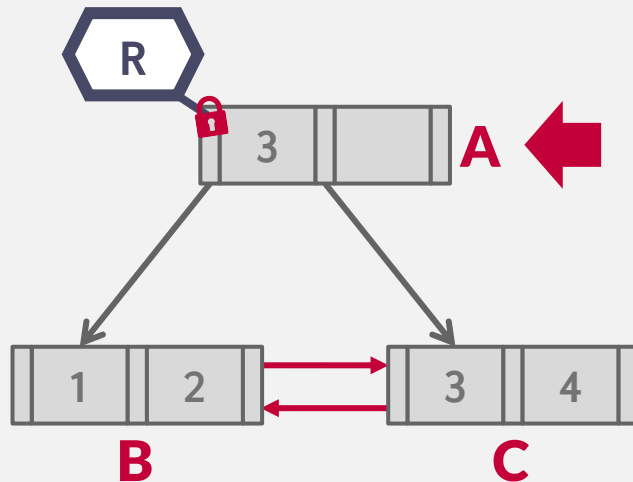
OBSERVATION

The threads in all the examples so far have acquired latches in a "top-down" manner.

- A thread can only acquire a latch from a node that is below its current node.
- If the desired latch is unavailable, the thread must wait until it becomes available.

But what if threads want to move from one leaf node to another leaf node?

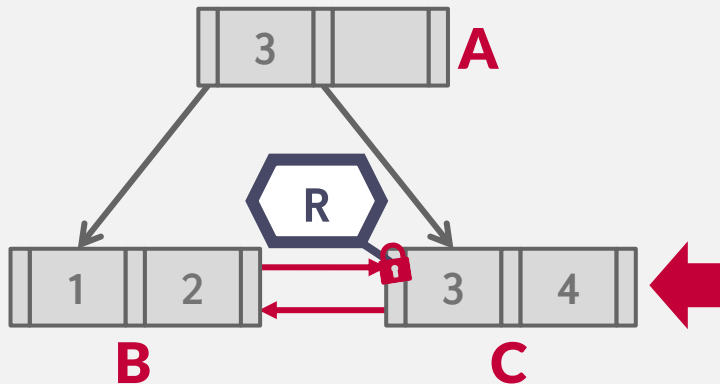
LEAF NODE SCAN EXAMPLE #1



T_1 : Find Keys < 4

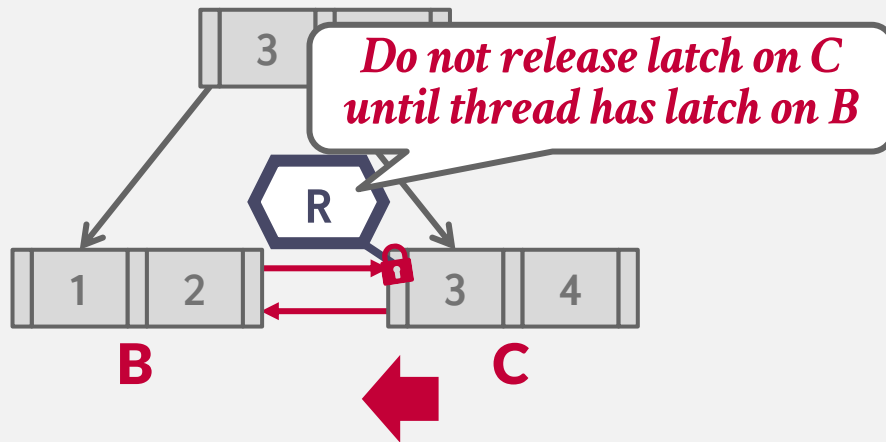
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



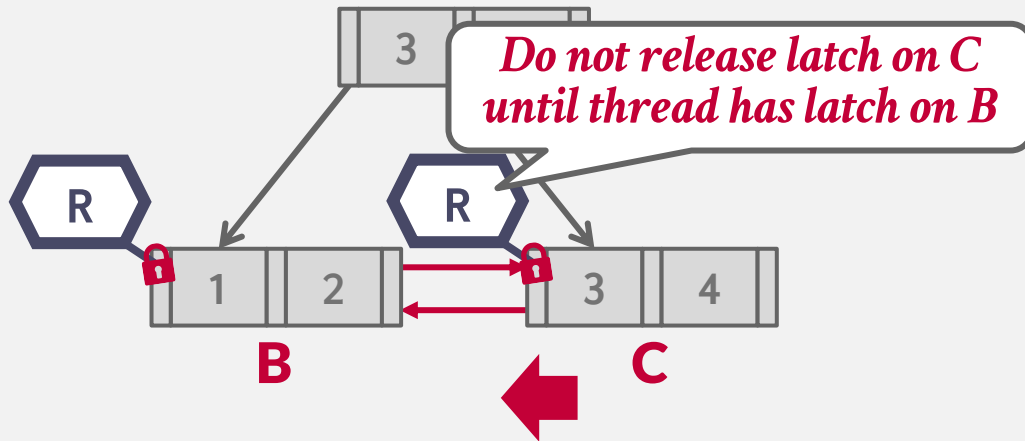
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



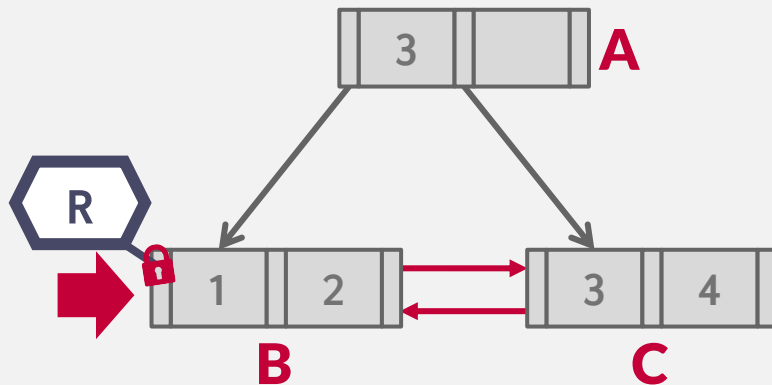
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4

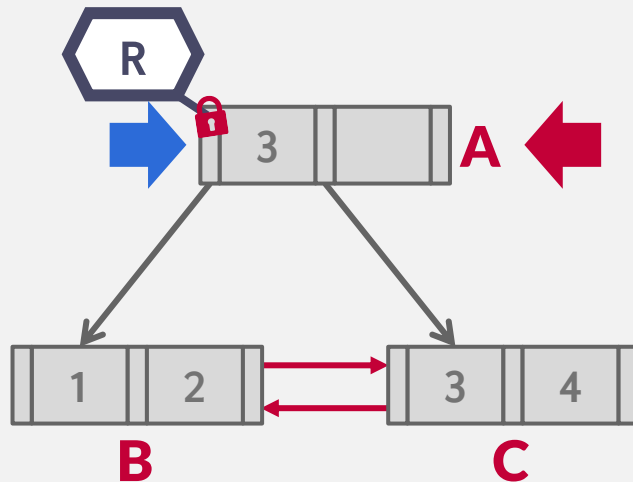


LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



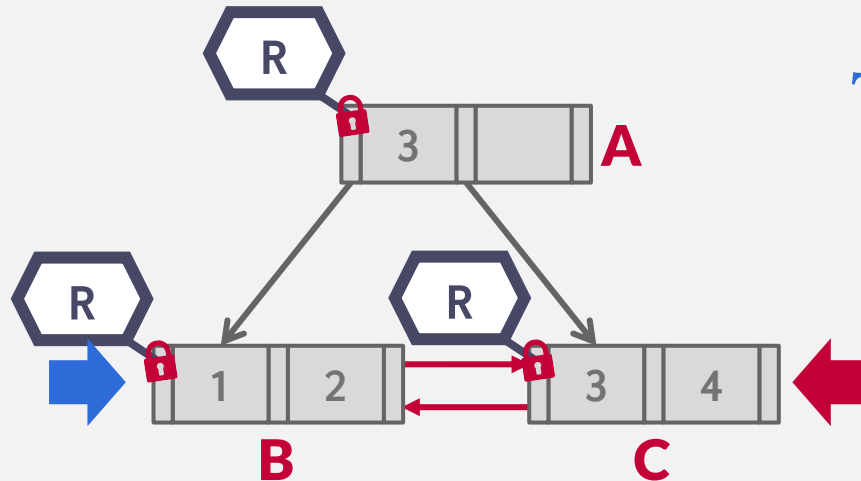
LEAF NODE SCAN EXAMPLE #2



T_1 : Find Keys < 4

T_2 : Find Keys > 1

LEAF NODE SCAN EXAMPLE #2



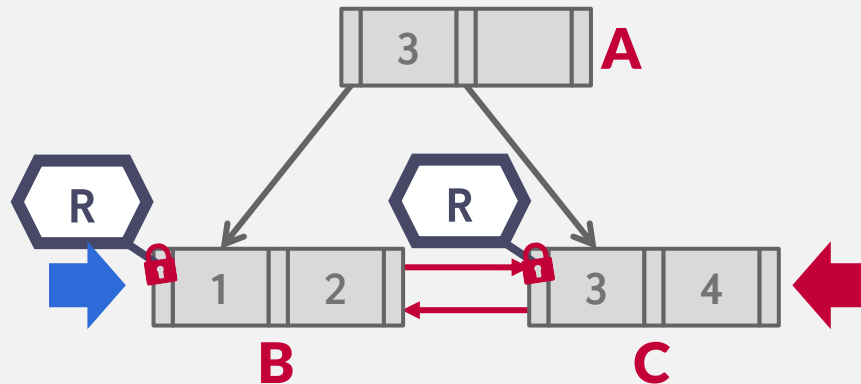
T_1 : Find Keys < 4

T_2 : Find Keys > 1

LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

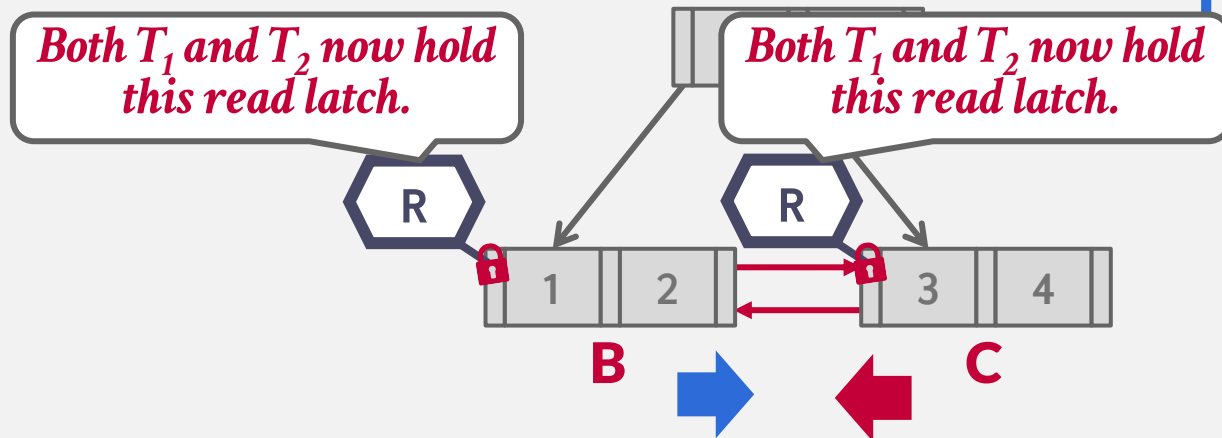
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

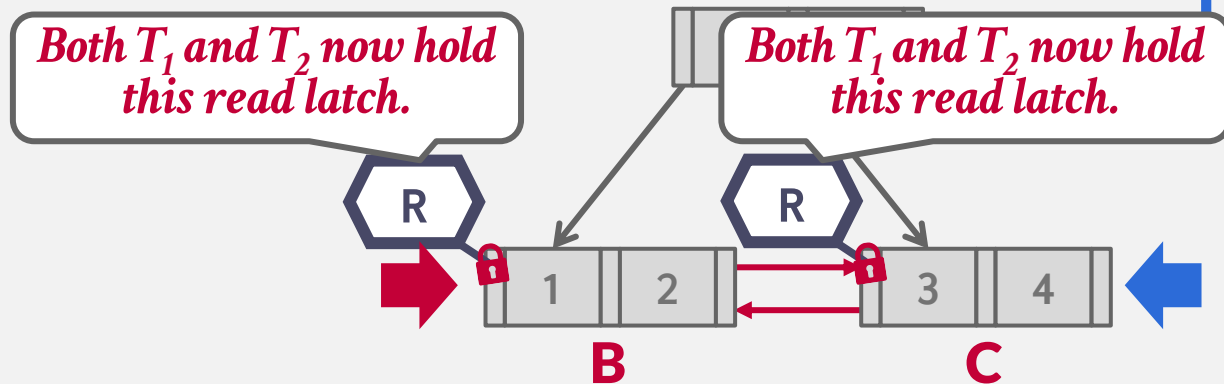
T_2 : Find Keys > 1



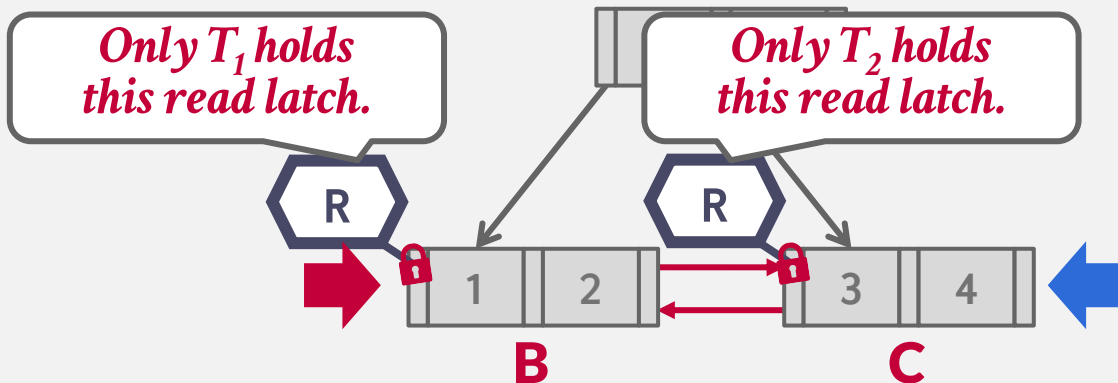
LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2



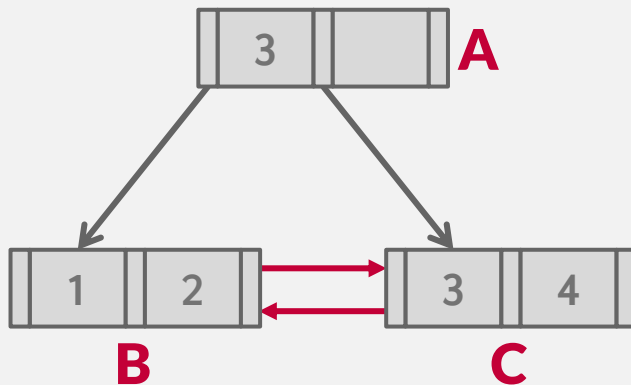
T_1 : Find Keys < 4

T_2 : Find Keys > 1

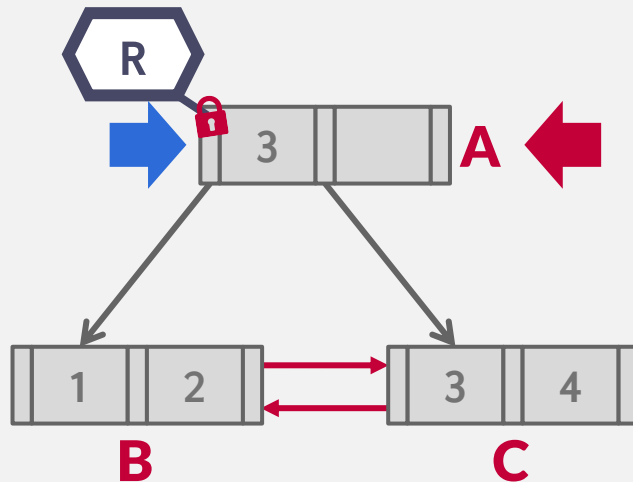
LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3



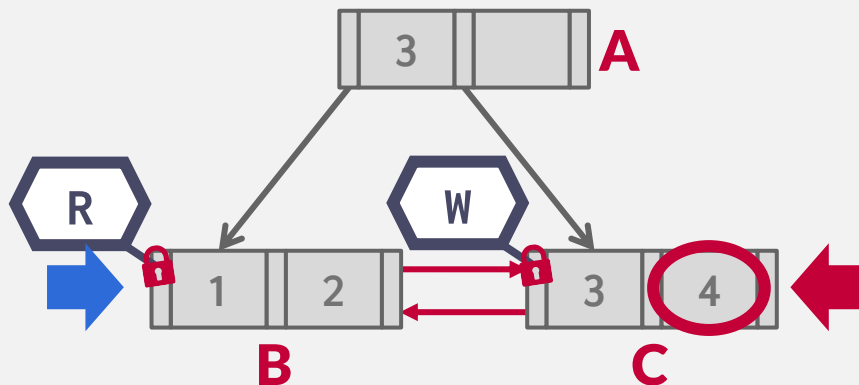
T_1 : Delete 4

T_2 : Find Keys > 1

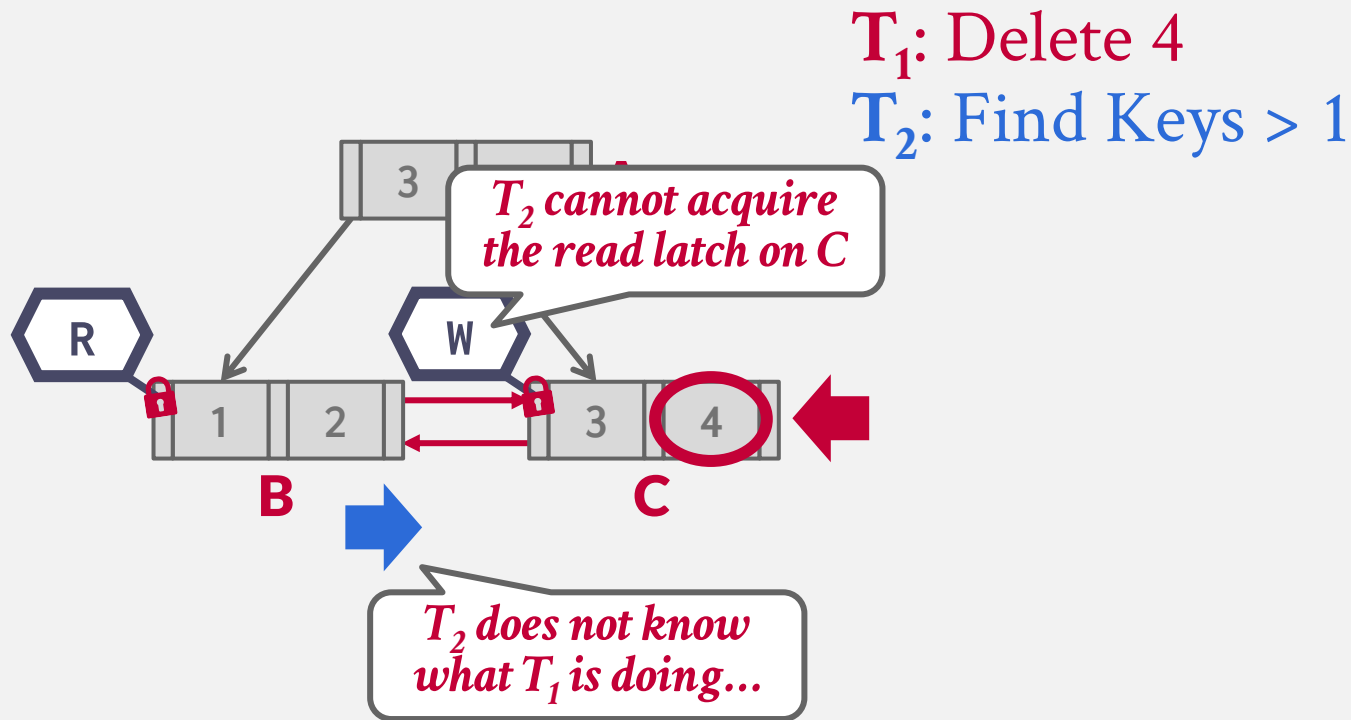
LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

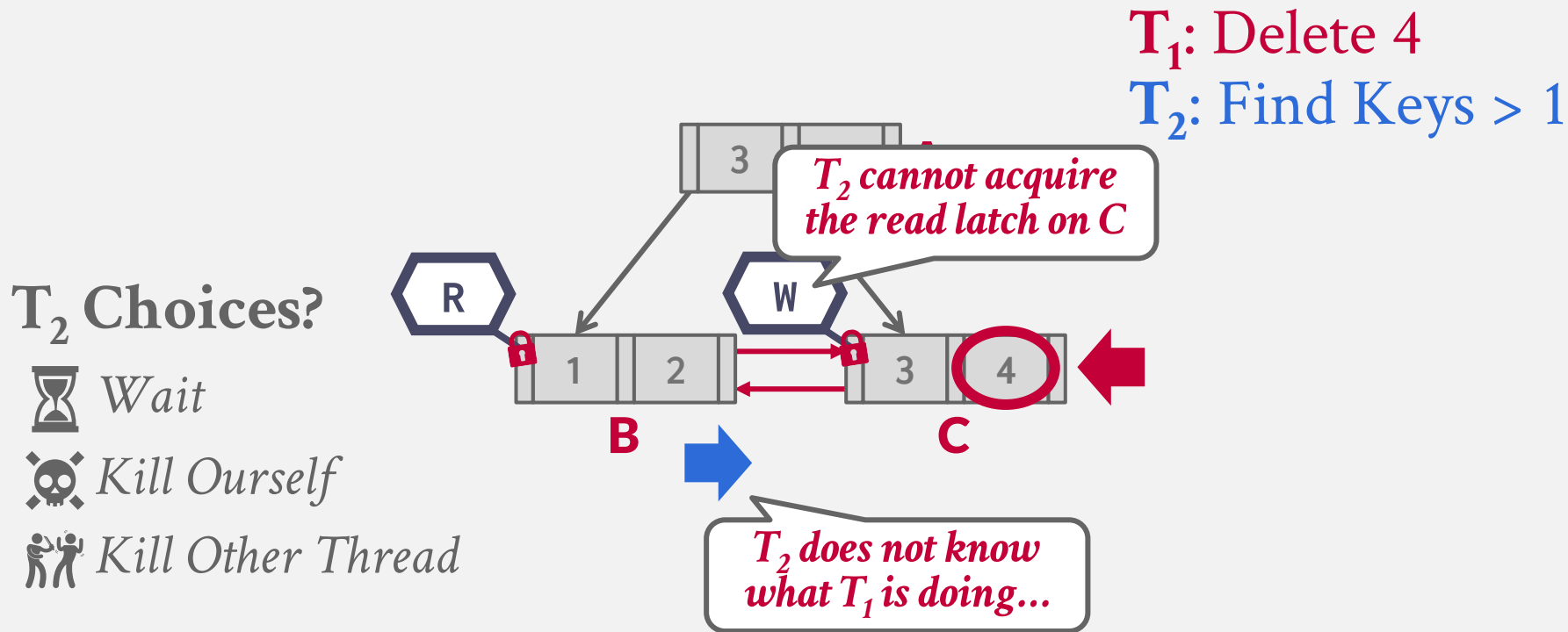
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3



LEAF NODE SCAN EXAMPLE #3



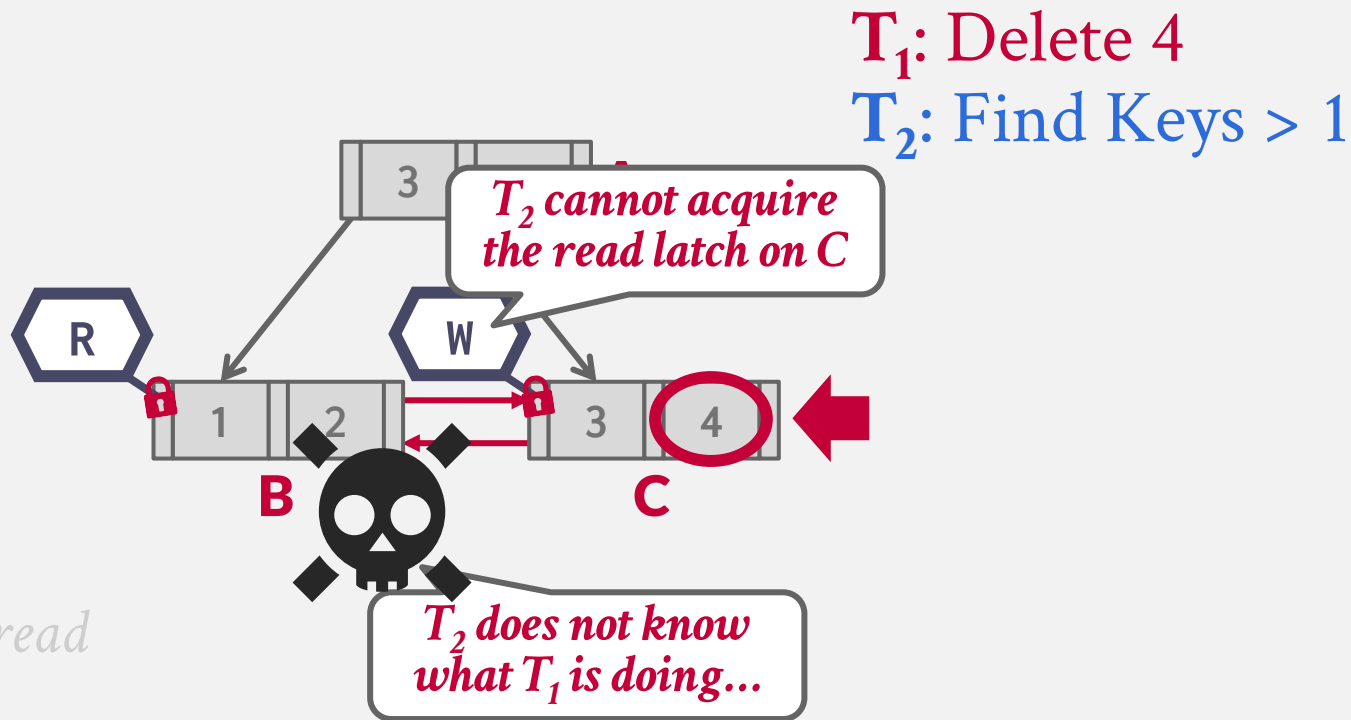
LEAF NODE SCAN EXAMPLE #3

T₂ Choices?

Wait

☠ Kill Ourselves

Kill Other Thread



LEAF NODE SCANS

Latches do not support deadlock detection or avoidance. The only way we can deal with this problem is through coding discipline.

The leaf node sibling latch acquisition protocol must support a “no-wait” mode.

The DBMS's data structures must cope with failed latch acquisitions.

CONCLUSION

Making a data structure thread-safe is notoriously difficult in practice.

We focused on B+Trees, but the same high-level techniques are applicable to other data structures.

NEXT CLASS

We are finally going to discuss how to execute some queries...

COMPARE-AND-SWAP

Atomic instruction that compares contents of a memory location **M** to a given value **V**

→ If values are equal, installs new given value **V'** in **M**

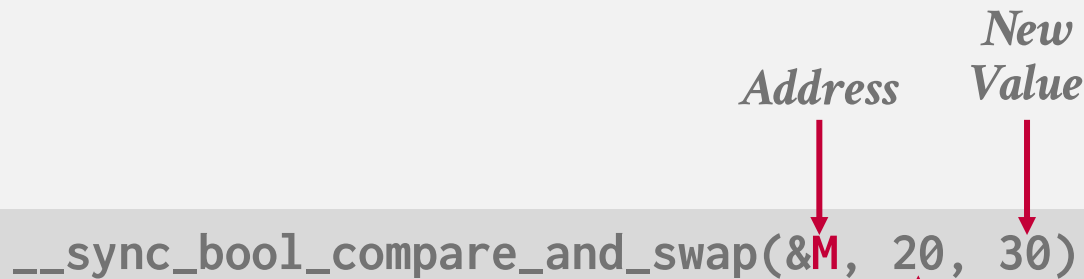
→ Otherwise, operation fails



Address *New Value*

```
__sync_bool_compare_and_swap(&M, 20, 30)
```

Compare Value



The function call is shown in a grey box. Red arrows point from the labels 'Address', 'New Value', and 'Compare Value' to the arguments &M, 20, and 30 respectively.

