



# Lecture 16: Index Concurrency Control

CREATING THE NEXT®

# Administrivia

---

- Project updates due on Nov 1 and Nov 3
- Assignment 3 and Sheet 3 due on Oct 27

# Today's Agenda

---

Recap

Latches Overview

Hash Table Latching

B+Tree Concurrency Control

Leaf Node Scans

*B*<sup>blink</sup>-Tree

Conclusion



# Recap







# Latches Overview











# Latch Implementations

---

- Blocking OS Mutex
- Test-and-Set Spin Latch
- Reader-Writer Latch



# Latch Implementations

- Approach 2: Test-and-Set Spin Latch (TAS)

- ▶ Very efficient (single instruction to latch/unlatch)
- ▶ Non-scalable, not cache friendly
- ▶ Example: `std::atomic<T>`
- ▶ Unlike OS mutex, spin latches do **not** suspend thread execution
- ▶ Atomic operations are faster if contention between threads is sufficiently **low**

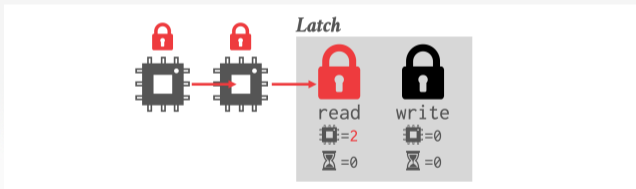
`std::atomic_flag latch; // atomic of boolean type (lock-free)`

```
while (latch.test_and_set(...)) {
  ~|// Retry? Yield? Abort?
}
```

# Latch Implementations

- Approach 3: Reader-Writer Latch**

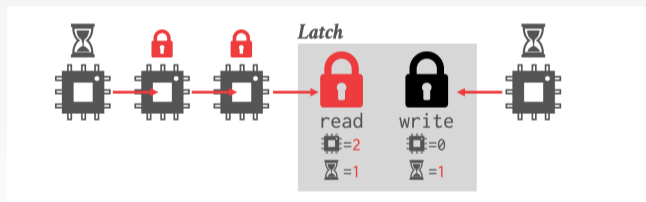
- ▶ Allows for concurrent readers
- ▶ Must manage read/write queues to avoid starvation
- ▶ Can be implemented on top of spinlocks



# Latch Implementations

- Approach 3: Reader-Writer Latch

- ▶ Allows for concurrent readers
- ▶ Must manage read/write queues to avoid starvation
- ▶ Can be implemented on top of spinlocks







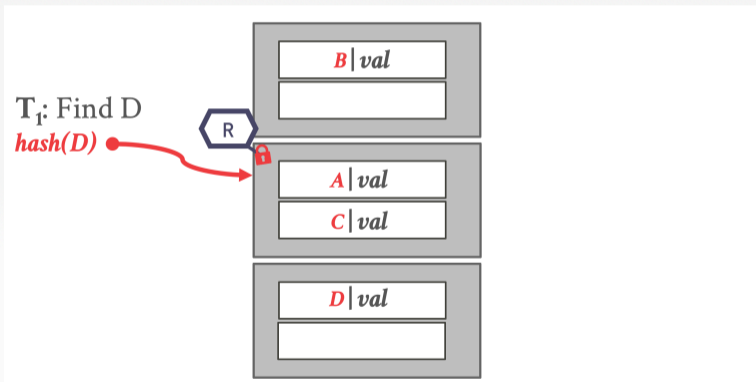


# Hash Table Latching

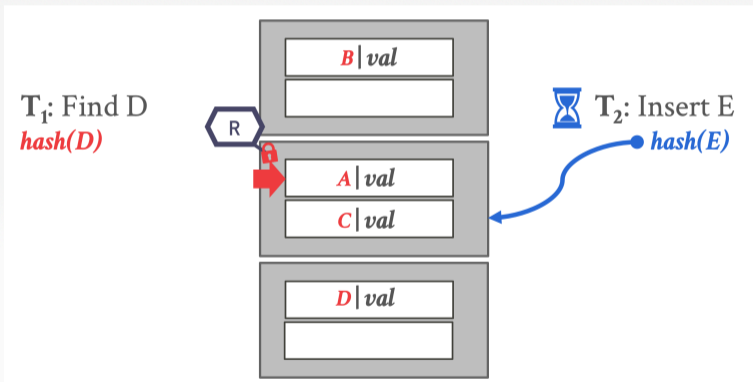
---

- **Approach 1: Page Latches**
  - ▶ Each page has its own reader-write 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.

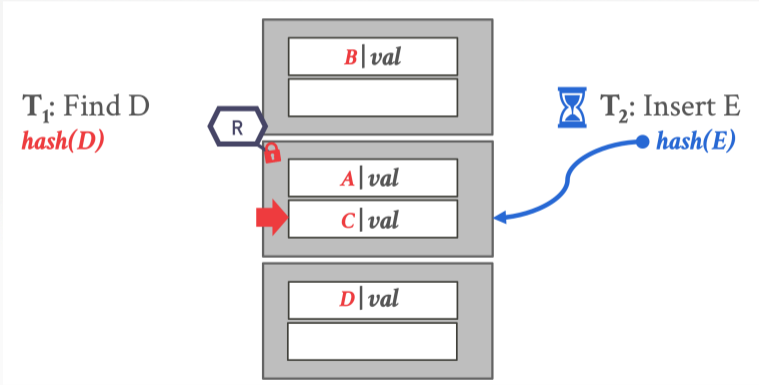
# Hash Table - Page Laches



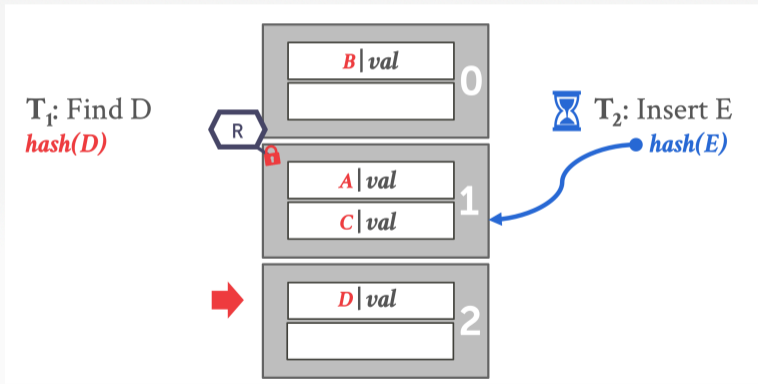
# Hash Table - Page Laches



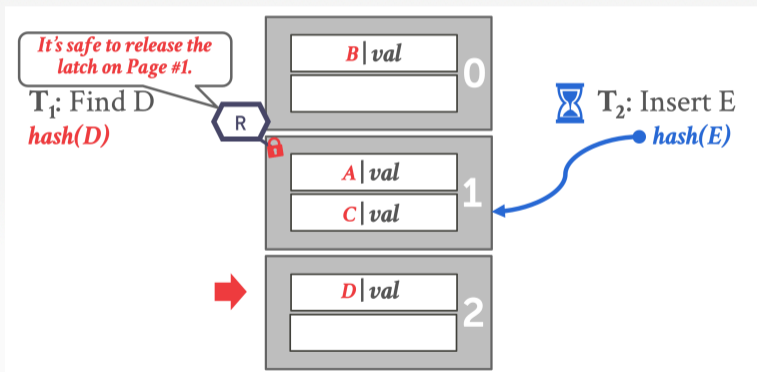
# Hash Table - Page Laches



# Hash Table - Page Laches

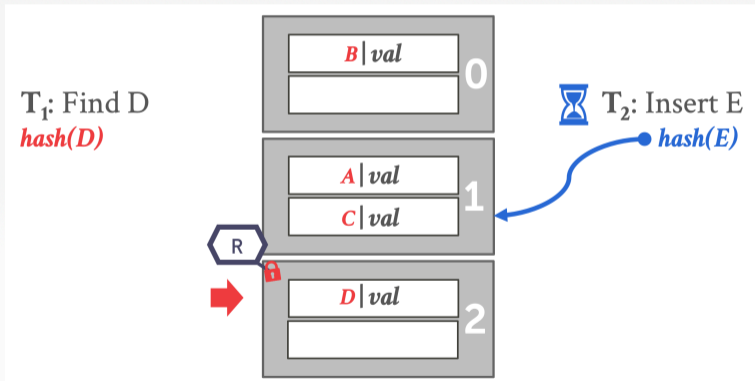


# Hash Table - Page Latches

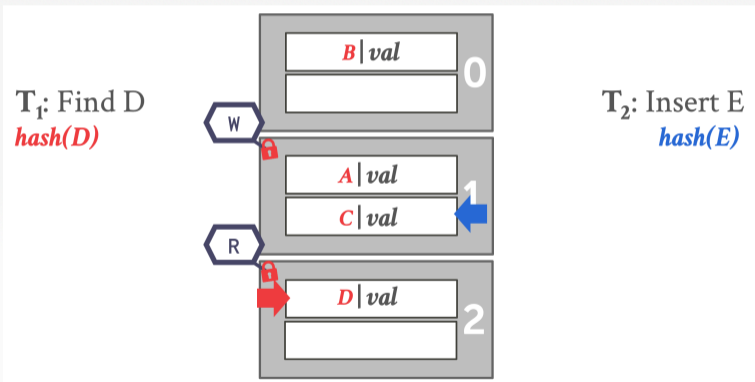




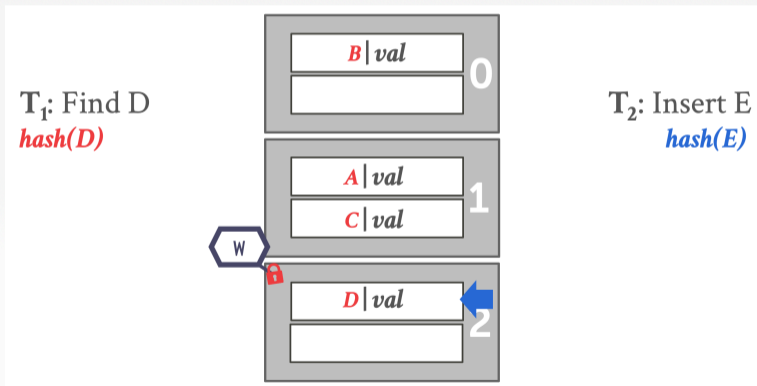
# Hash Table - Page Laches



# Hash Table - Page Laches



# Hash Table - Page Laches

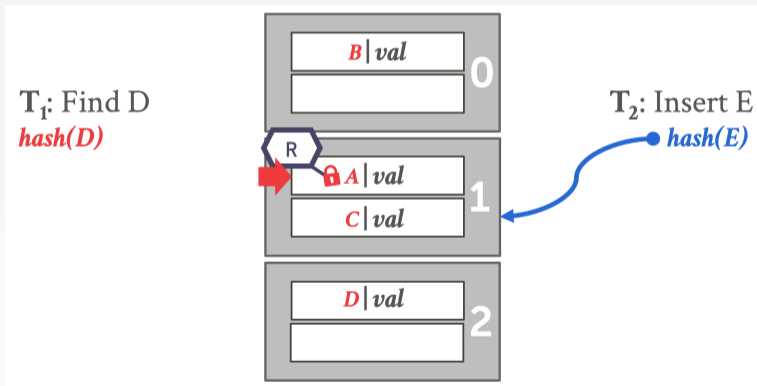




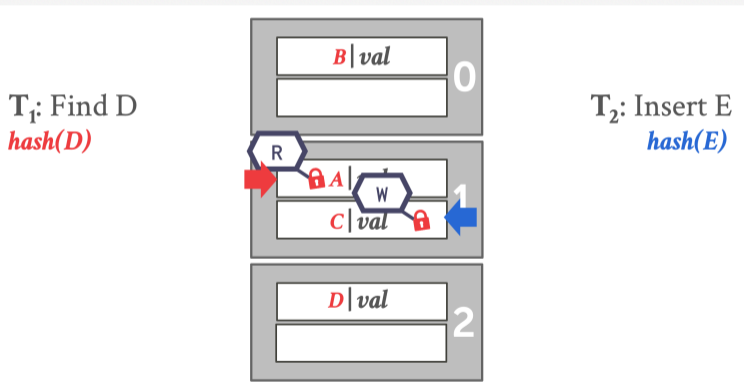




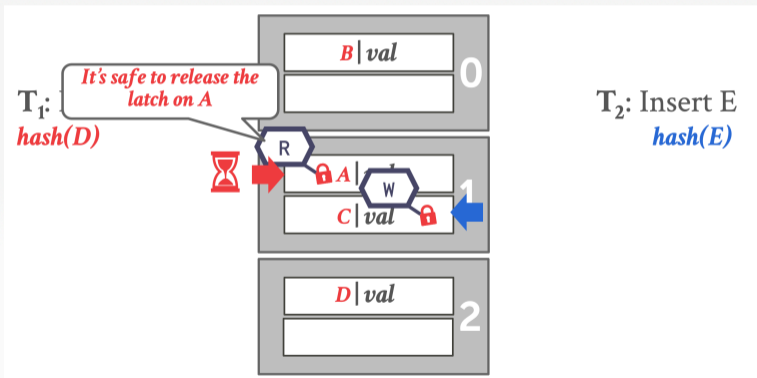
# Hash Table - Slot Laches



# Hash Table - Slot Laches



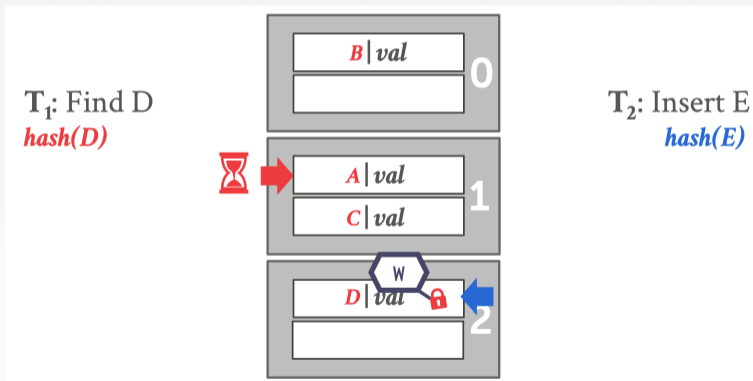
# Hash Table - Slot Laches







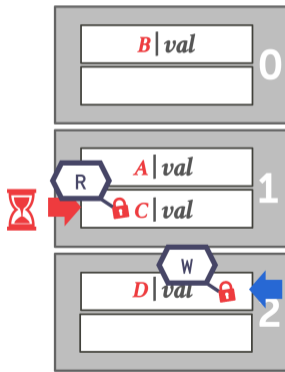
# Hash Table - Slot Laches





# Hash Table - Slot Latches

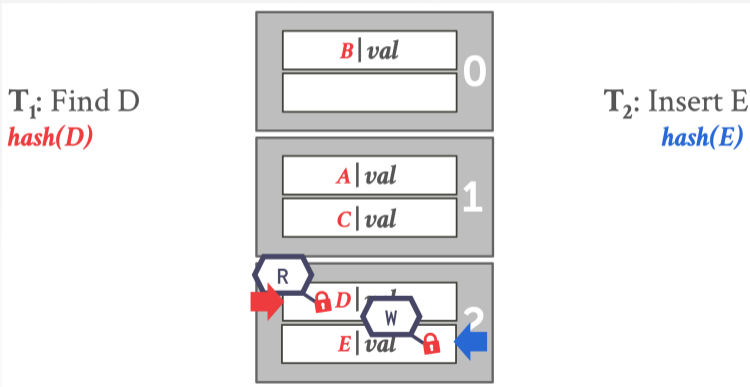
$T_1$ : Find D  
*hash(D)*



$T_2$ : Insert E  
*hash(E)*



# Hash Table - Slot Laches





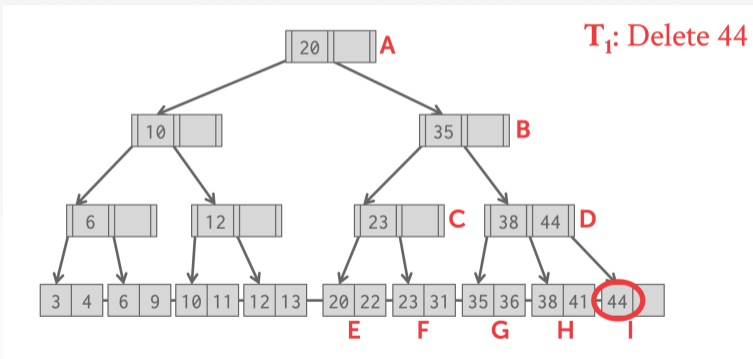
# B Tree Concurrency Control

# B+Tree Concurrency Control

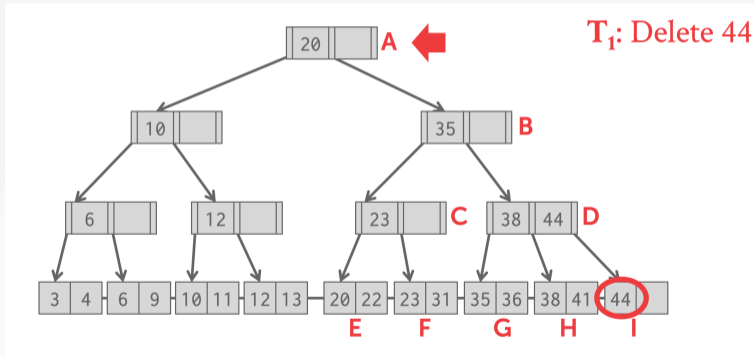
---

- We want to allow multiple threads to read and update a B+Tree at the same time.
- We need to handle 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 Concurrency Control: Example

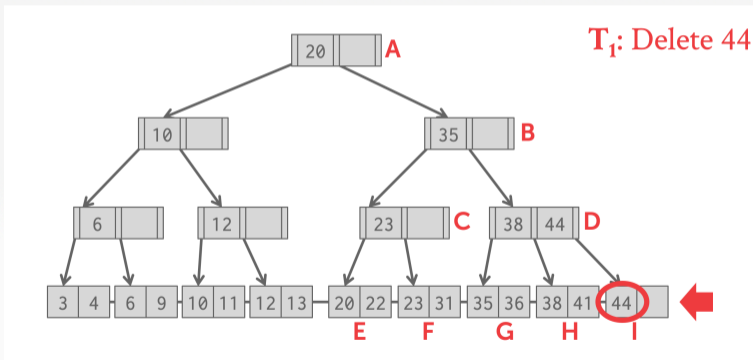


# B+Tree Concurrency Control: Example

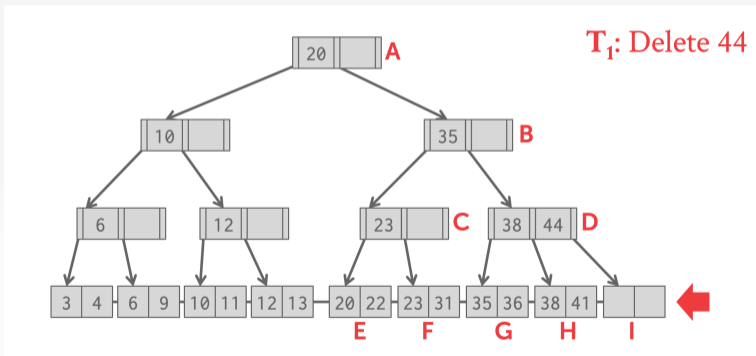




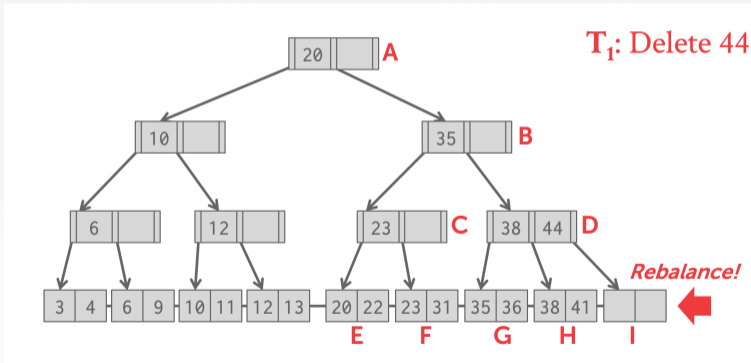
# B+Tree Concurrency Control: Example



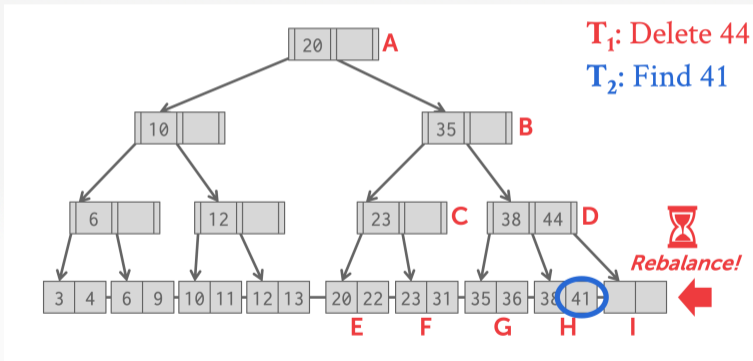
# B+Tree Concurrency Control: Example



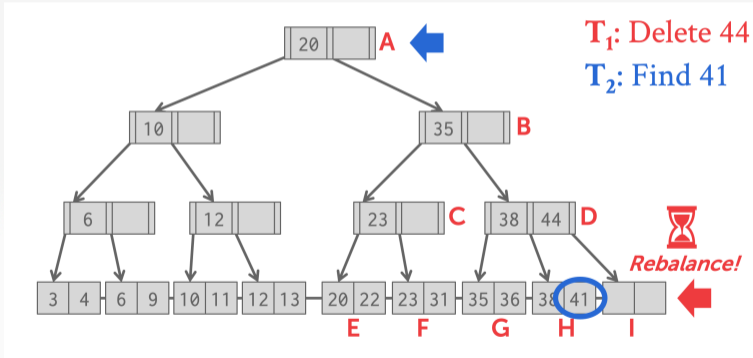
# B+Tree Concurrency Control: Example



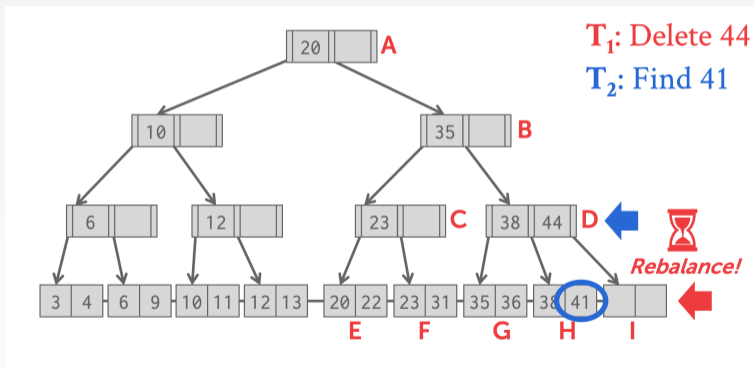
# B+Tree Concurrency Control: Example



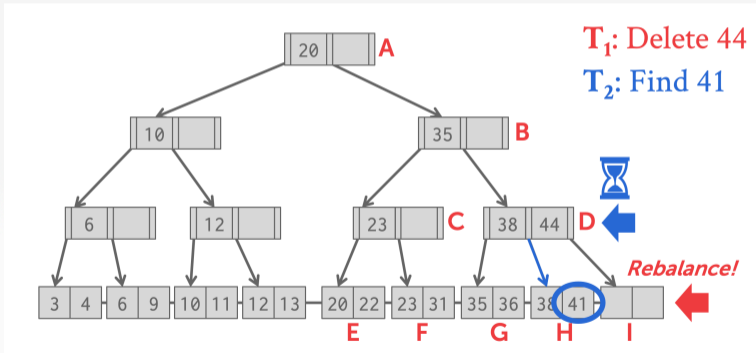
# B+Tree Concurrency Control: Example



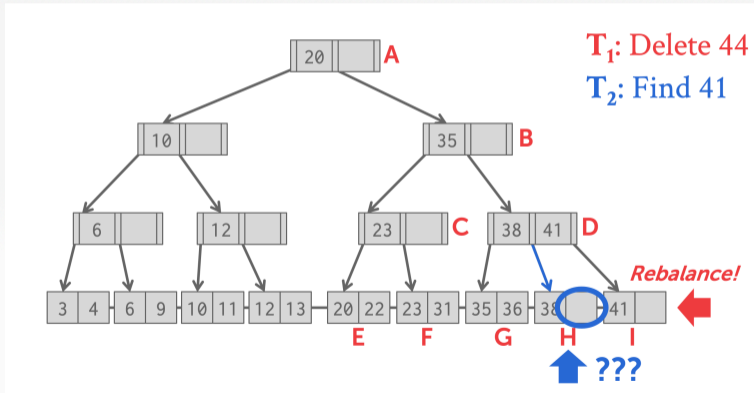
# B+Tree Concurrency Control: Example



# B+Tree Concurrency Control: Example



# B+Tree Concurrency Control: Example





# Latch Crabbing/Coupling

---

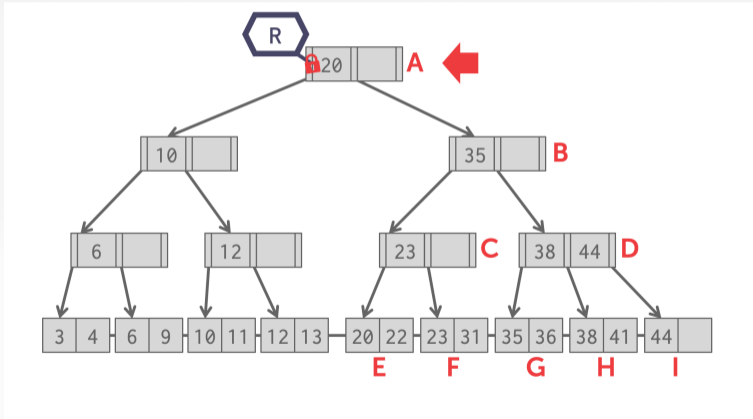
- Protocol to allow multiple threads to access/modify B+Tree at the same time.
- Basic Idea:
  - ▶ Get latch for parent.
  - ▶ Get latch for child
  - ▶ Release latch for parent if “safe”.
- A **safe node** is one that will **not split or merge** when updated.
  - ▶ Not full (on insertion)
  - ▶ More than half-full (on deletion)

# Latch Crabbing/Coupling

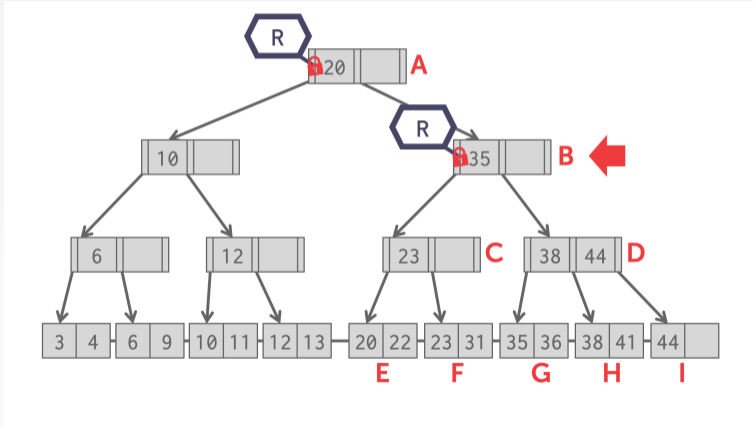
---

- **Find**: Start at root and go down; repeatedly,
  - ▶ Acquire **R** latch on child
  - ▶ Then unlatch parent
- **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.

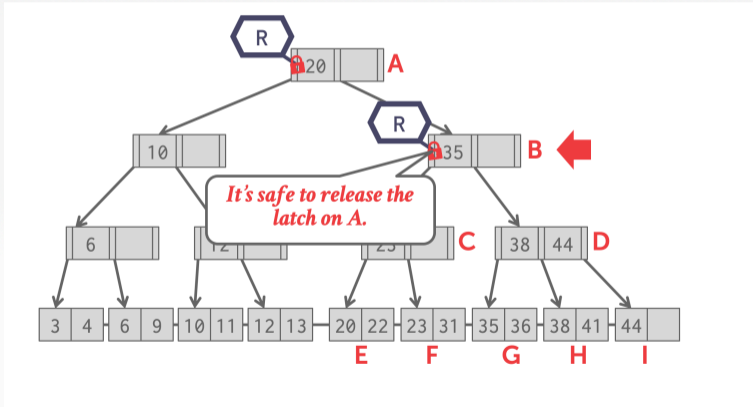
# Example 1 - Find 38



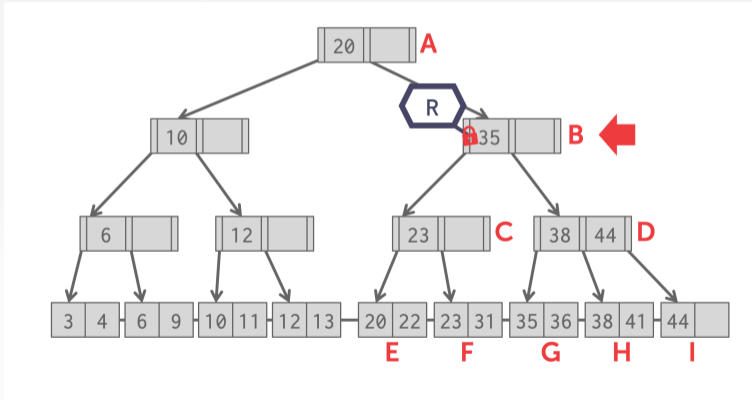
# Example 1 - Find 38



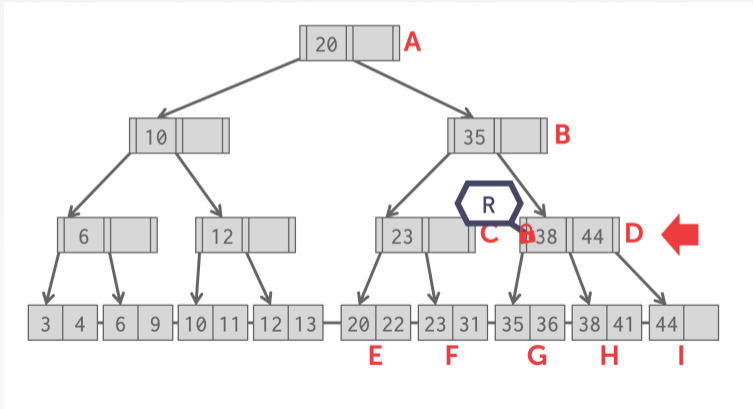
# Example 1 - Find 38



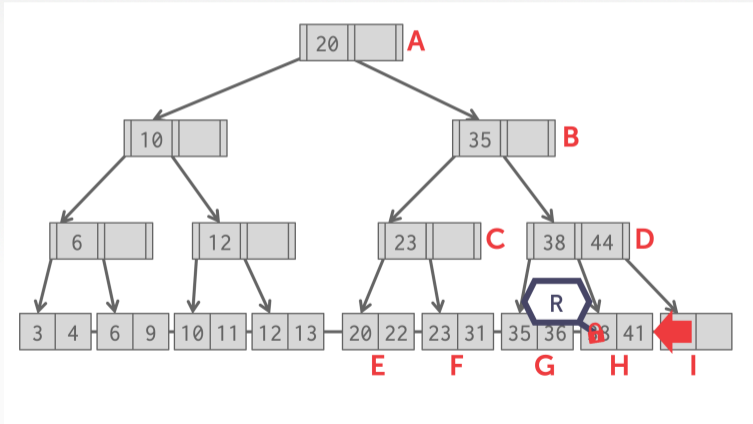
# Example 1 - Find 38



# Example 1 - Find 38

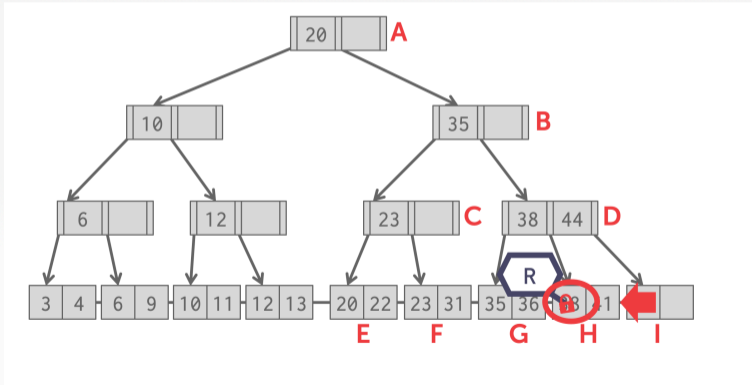


# Example 1 - Find 38

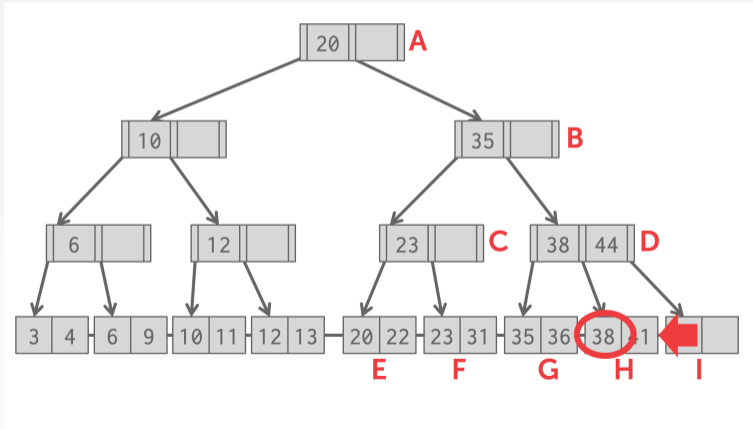




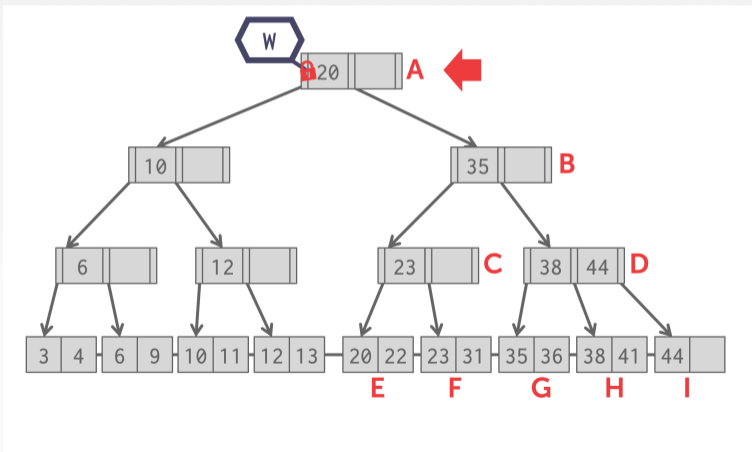
# Example 1 - Find 38



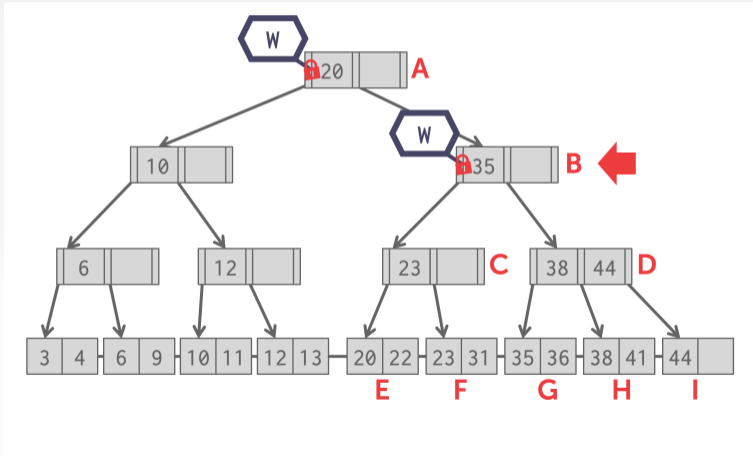
# Example 1 - Find 38



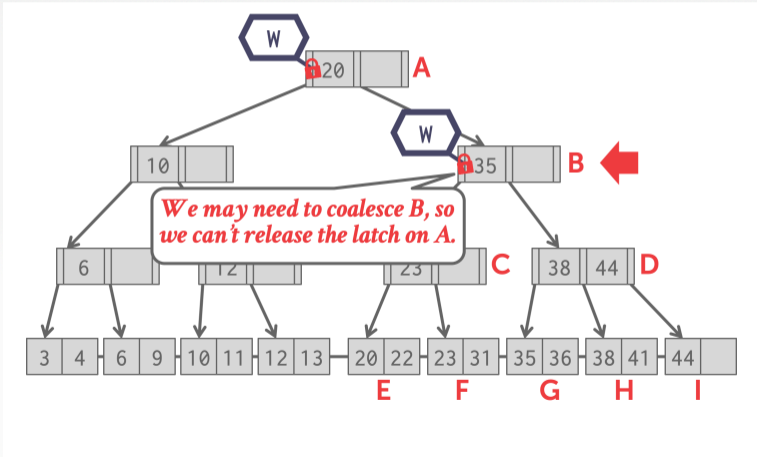
# Example 2 - Delete 38



# Example 2 - Delete 38

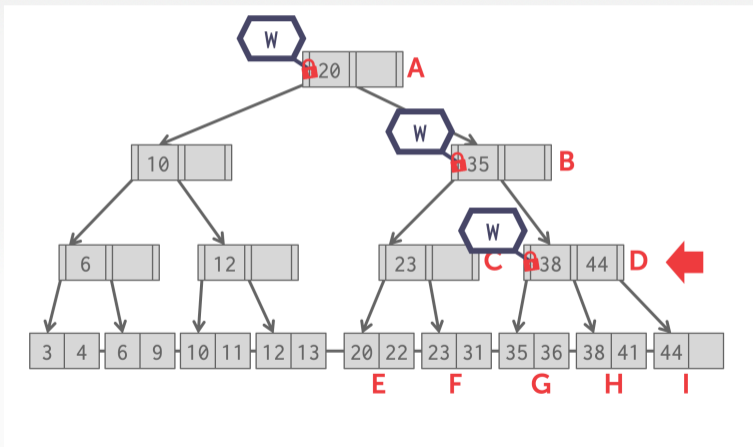


# Example 2 - Delete 38

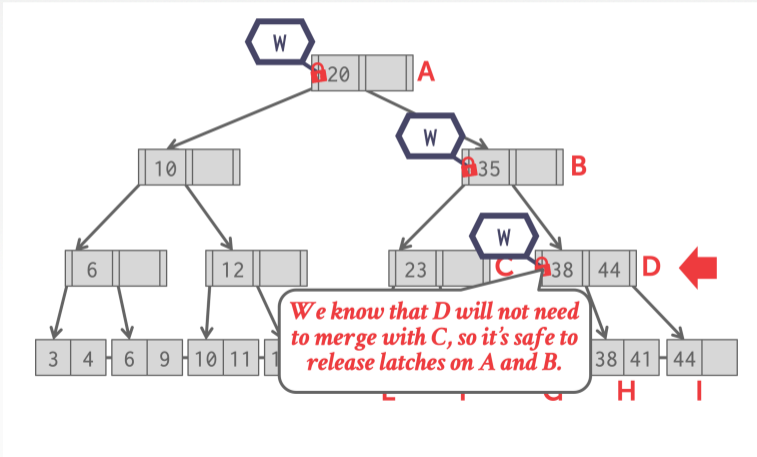




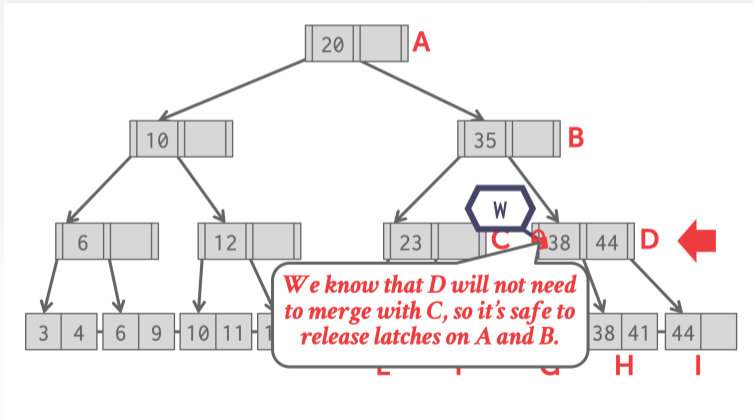
## Example 2 - Delete 38



# Example 2 - Delete 38

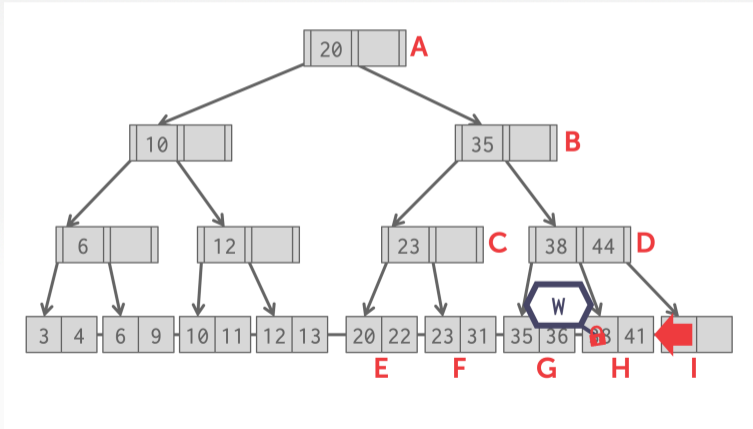


# Example 2 - Delete 38

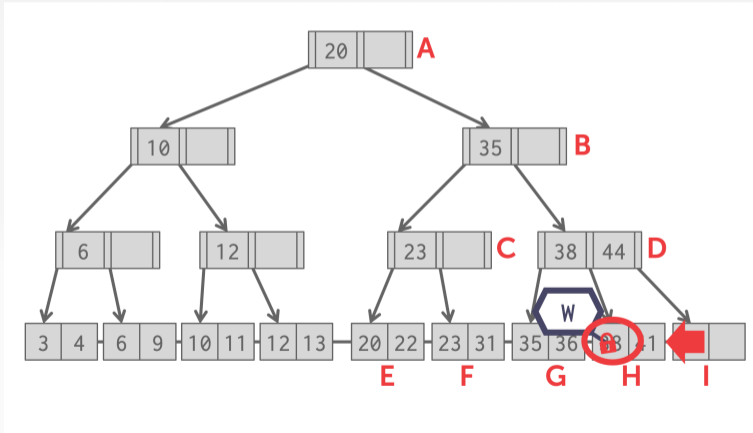




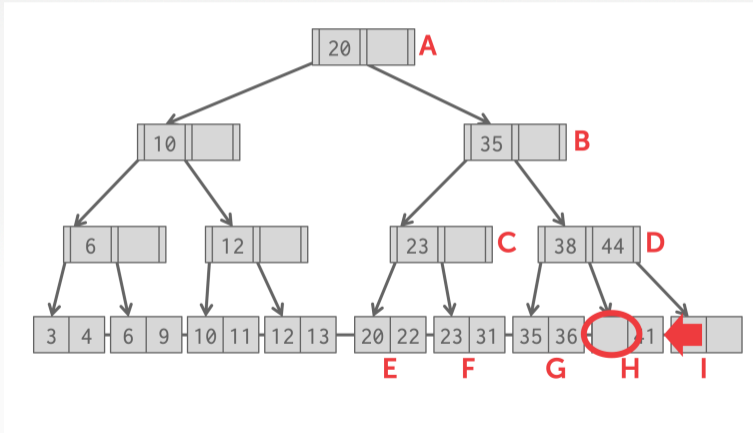
# Example 2 - Delete 38



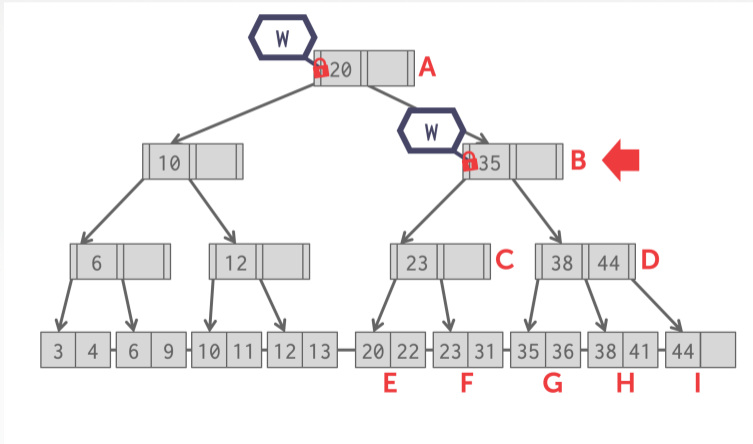
# Example 2 - Delete 38



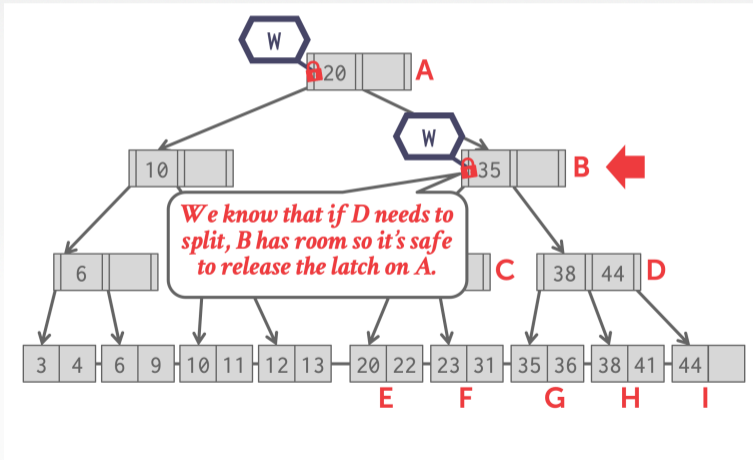
# Example 2 - Delete 38



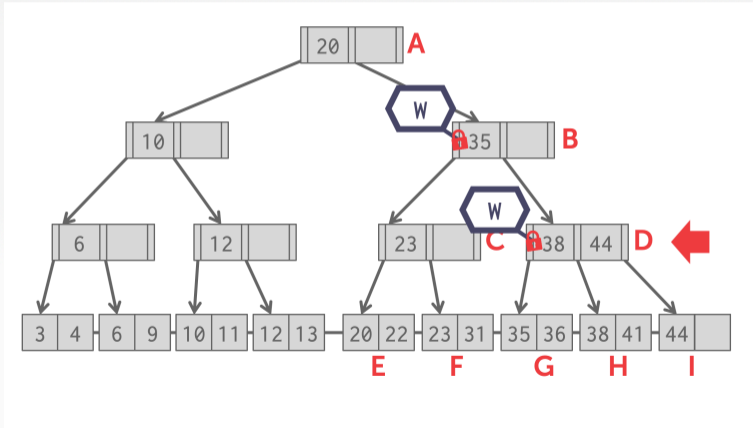
# Example 3 - Insert 45



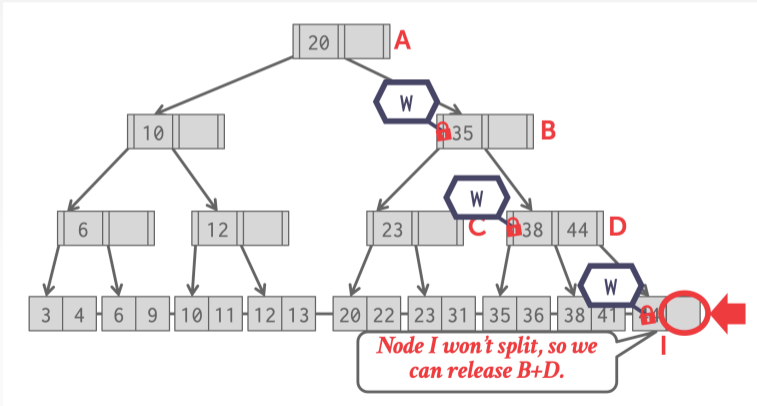
## Example 3 - Insert 45



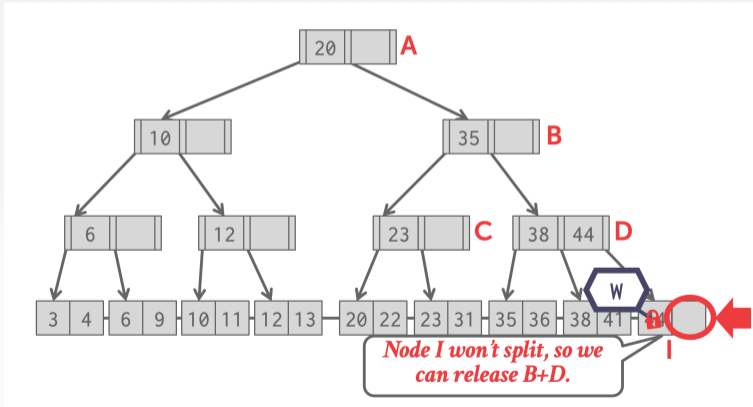
# Example 3 - Insert 45



# Example 3 - Insert 45

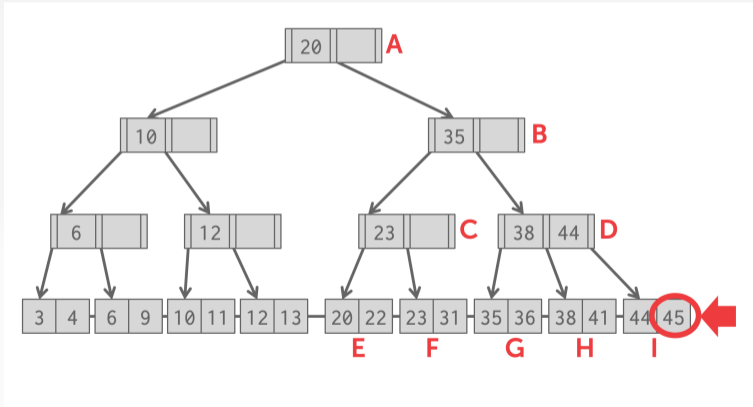


# Example 3 - Insert 45

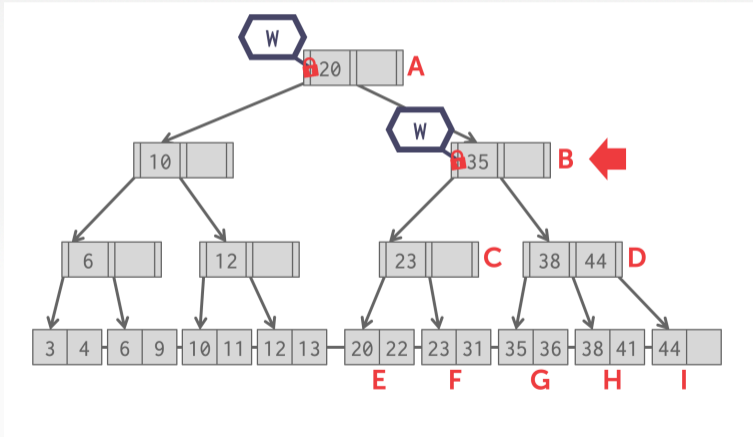




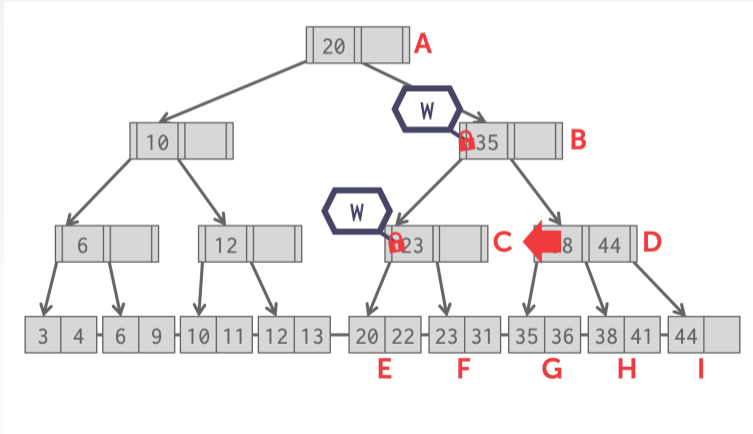
# Example 3 - Insert 45



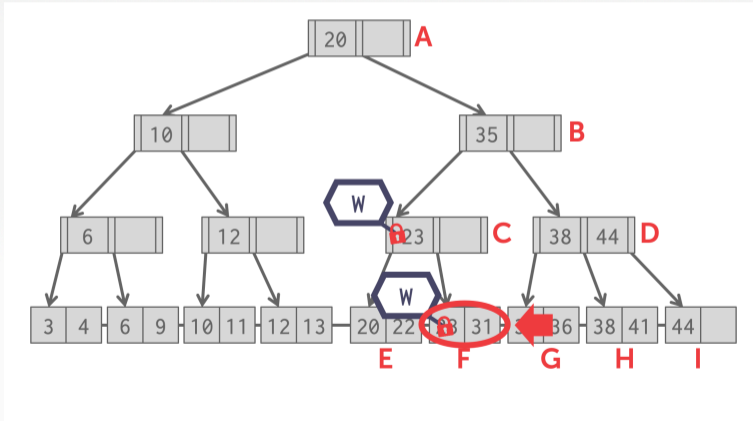
# Example 4 - Insert 25



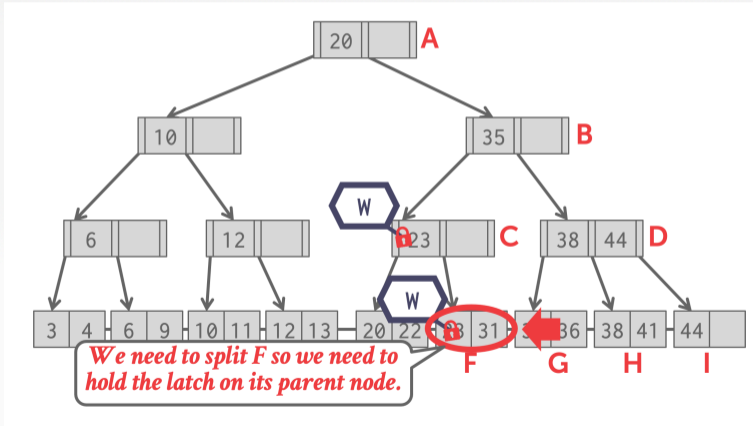
# Example 4 - Insert 25



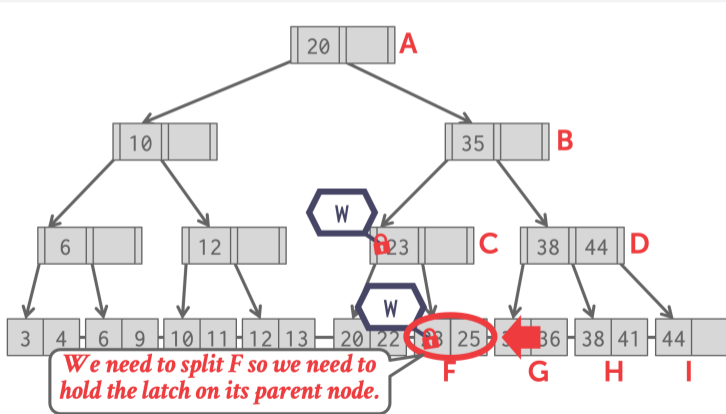
# Example 4 - Insert 25



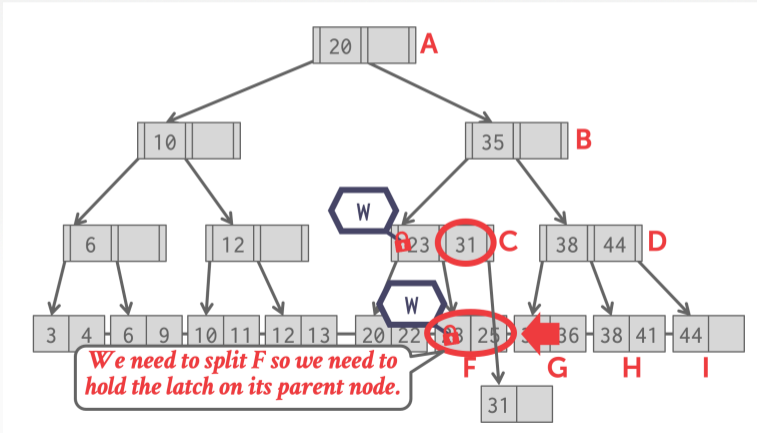
# Example 4 - Insert 25



# Example 4 - Insert 25



# Example 4 - Insert 25





## 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.
- Can we do better?

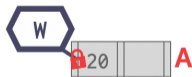
Delete 38



Insert 45



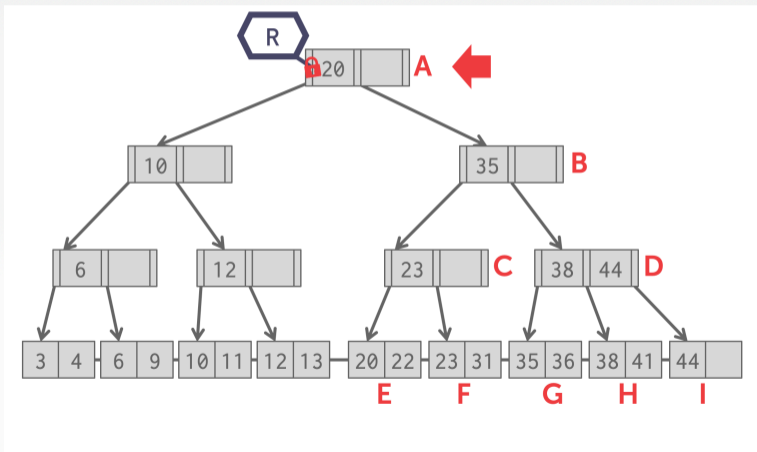
Insert 25



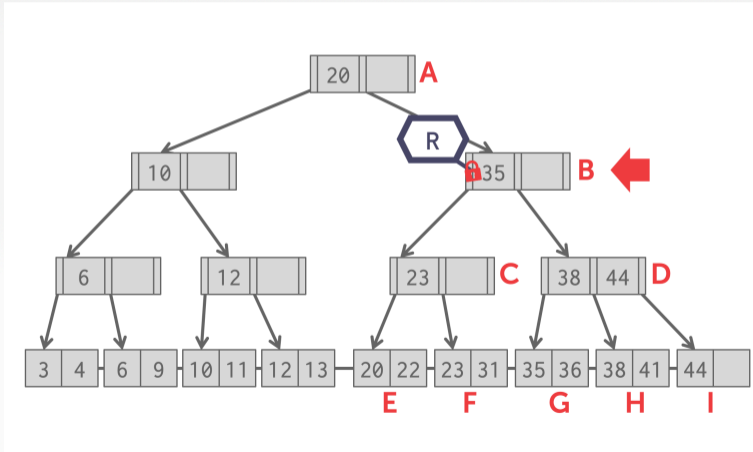




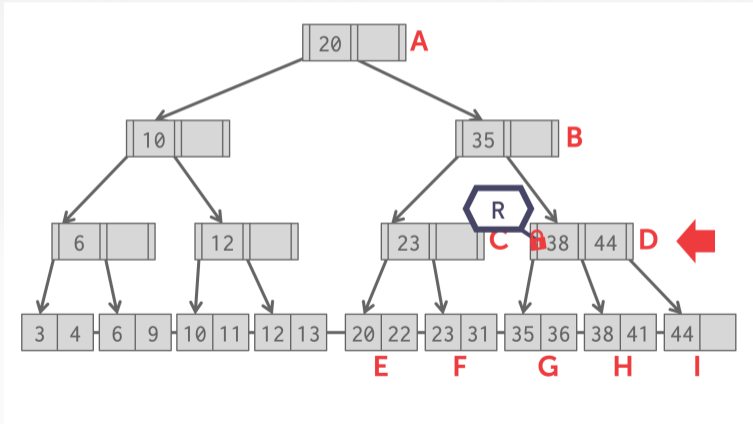
## Example 2 - Delete 38



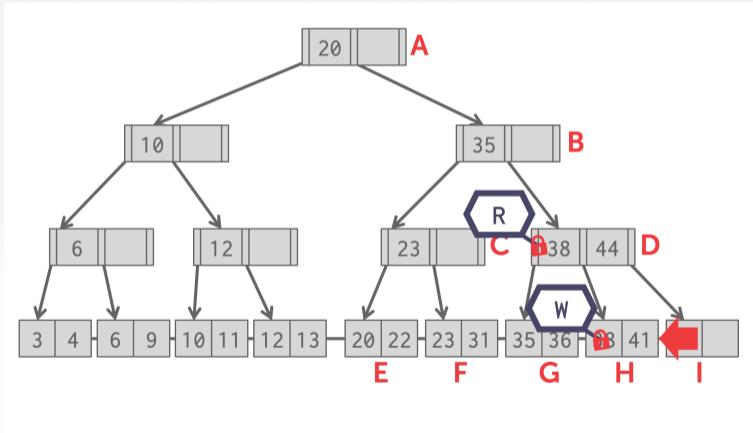
# Example 2 - Delete 38



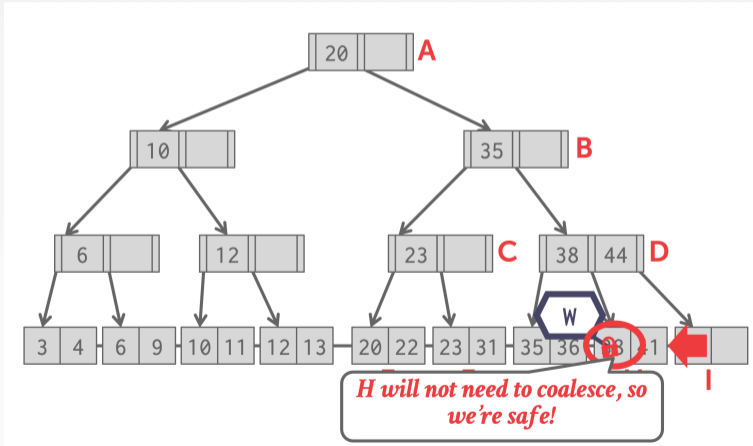
# Example 2 - Delete 38



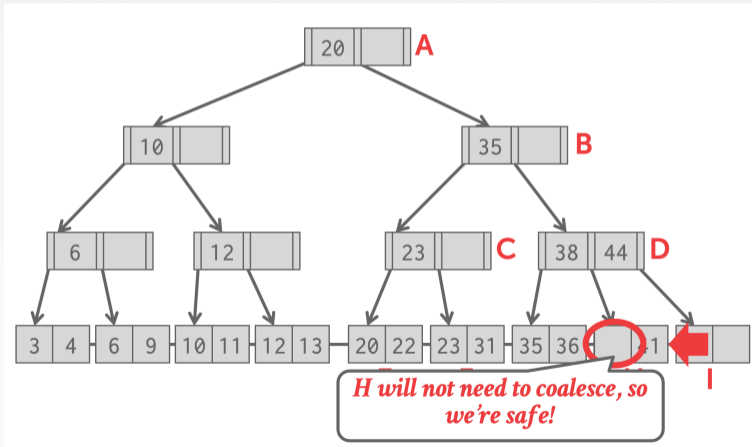
# Example 2 - Delete 38



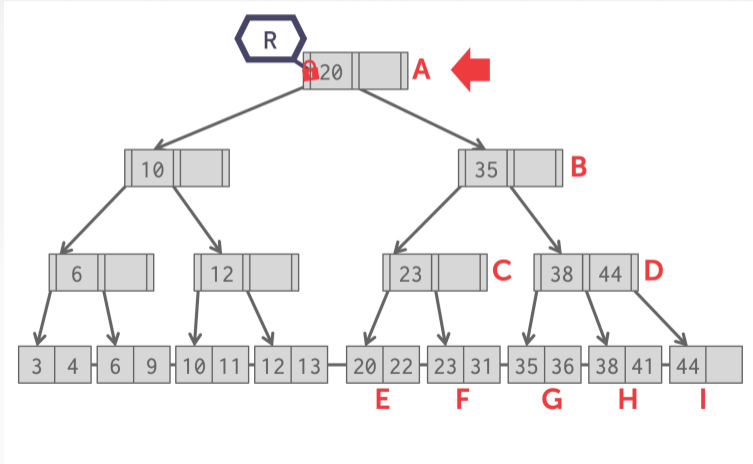
# Example 2 - Delete 38



# Example 4 - Insert 25

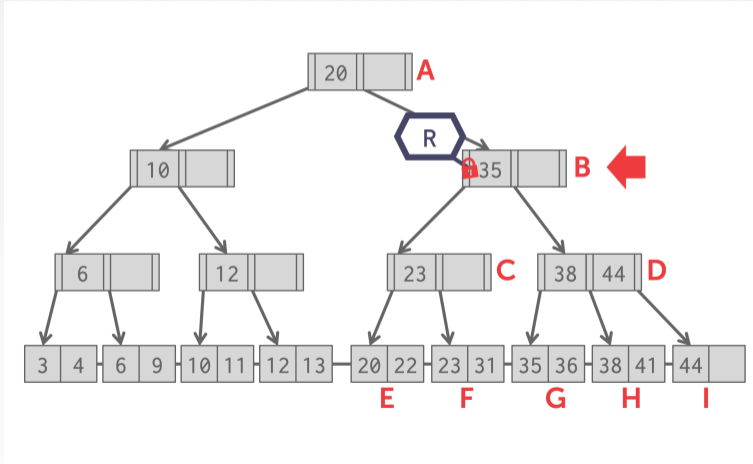


# Example 4 - Insert 25

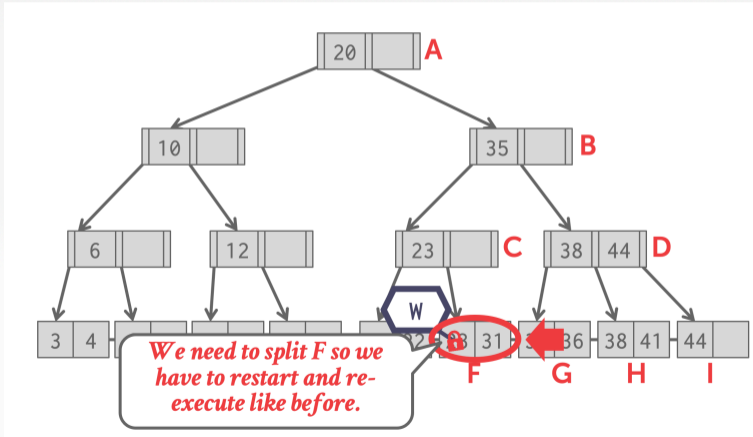




# Example 4 - Insert 25



# Example 4 - Insert 25







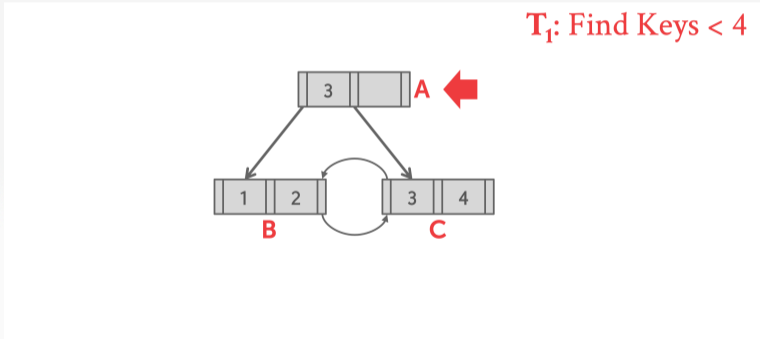
# Leaf Node Scans

## Observation

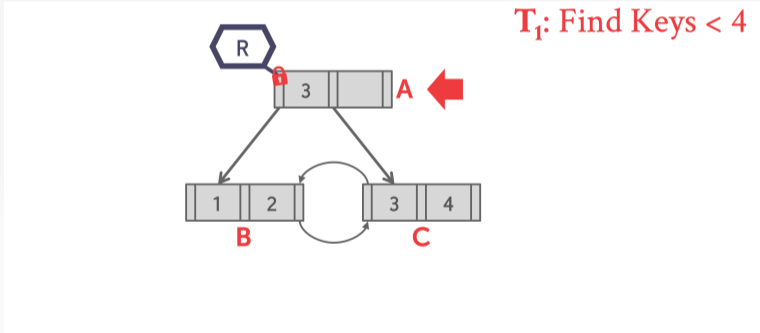
---

- 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 we want to move from one leaf node to another leaf node?
- Leaf nodes can include **hint keys** to approximate the next key at your sibling.

# Leaf Node Scan - Example 1

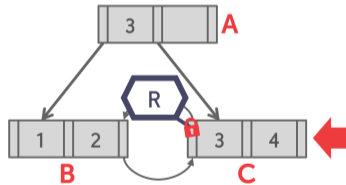


# Leaf Node Scan - Example 1



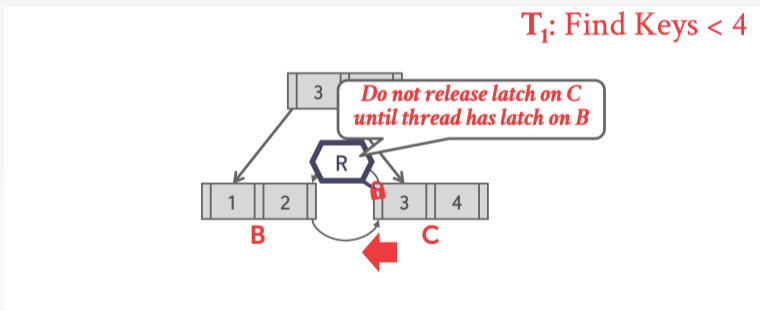
# Leaf Node Scan - Example 1

$T_1$ : Find Keys < 4

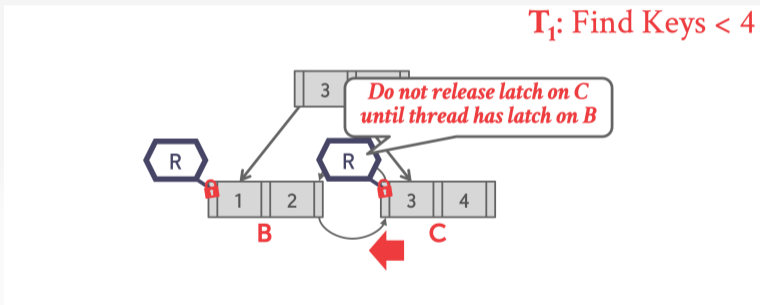




# Leaf Node Scan - Example 1

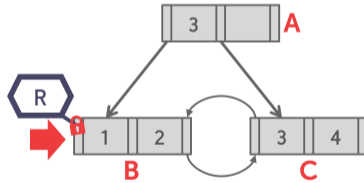


# Leaf Node Scan - Example 1

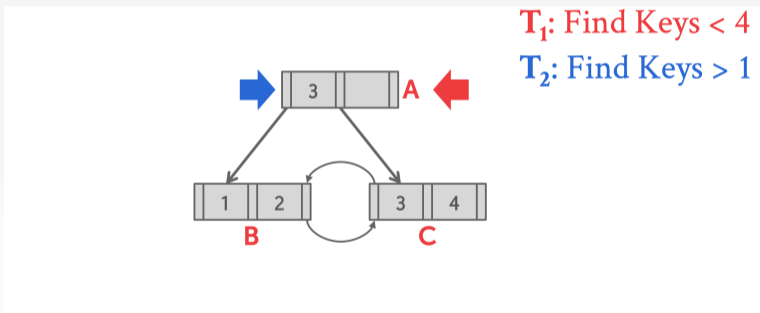


# Leaf Node Scan - Example 1

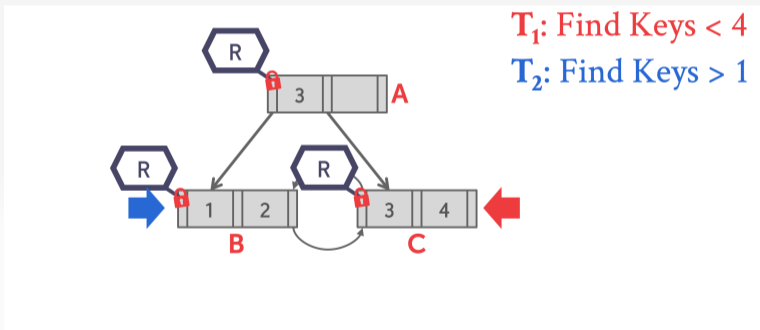
$T_1$ : Find Keys < 4



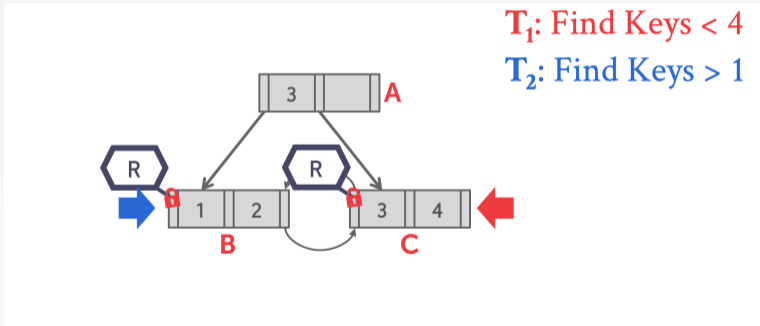
# Leaf Node Scan - Example 2



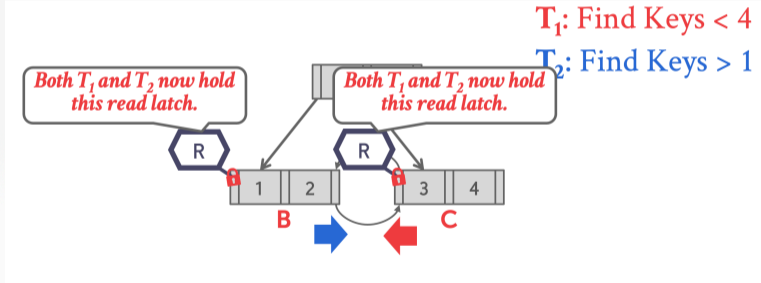
# Leaf Node Scan - Example 2



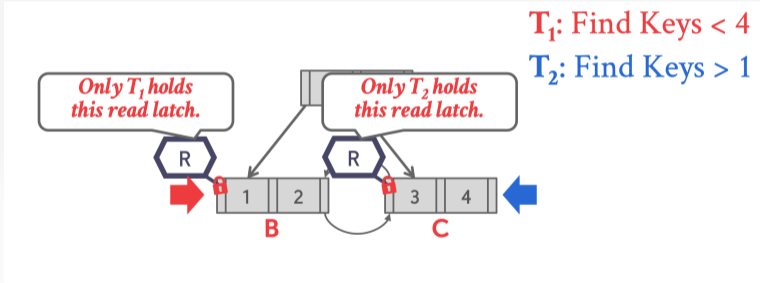
# Leaf Node Scan - Example 2



# Leaf Node Scan - Example 2

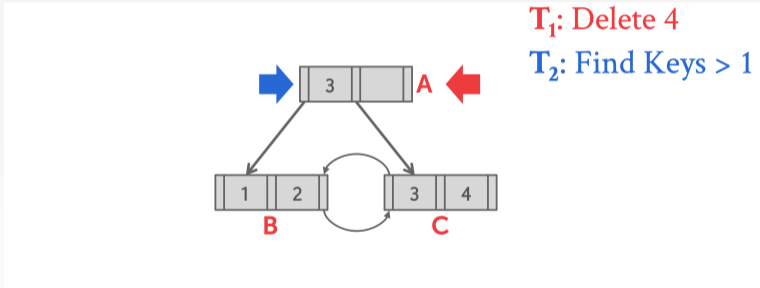


# Leaf Node Scan - Example 2

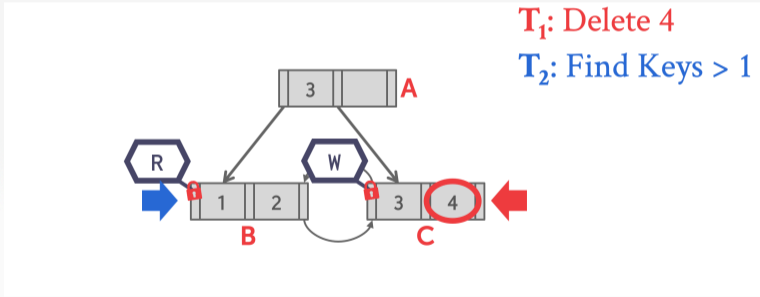




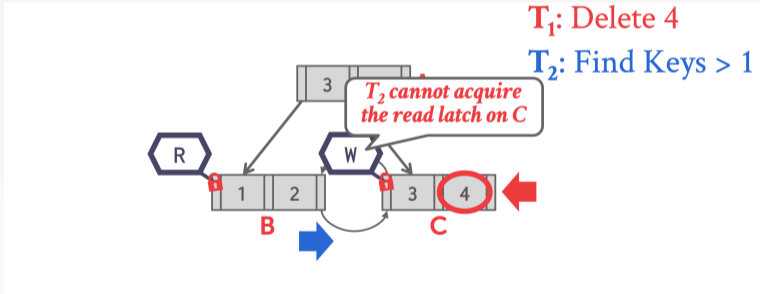
# Leaf Node Scan - Example 3



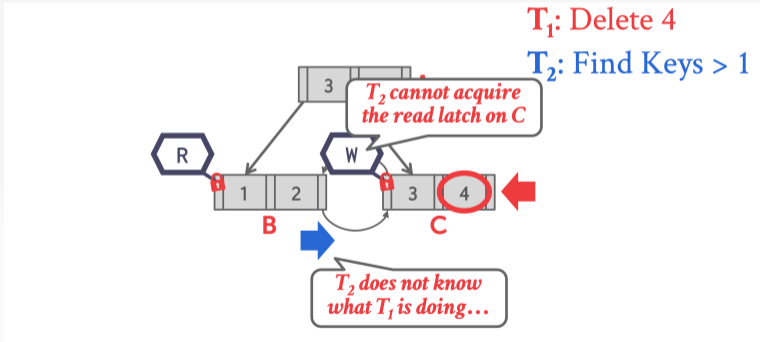
# Leaf Node Scan - Example 3



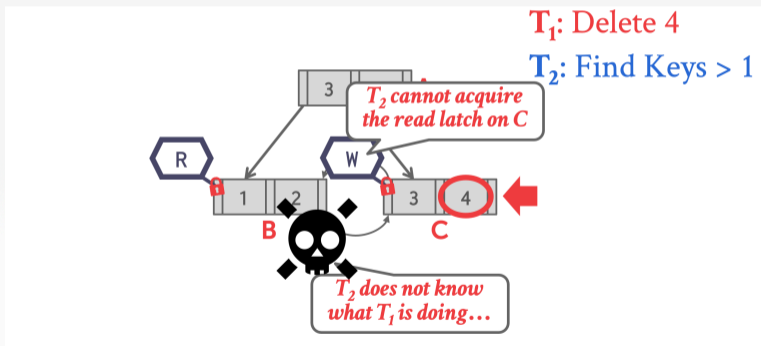
# Leaf Node Scan - Example 3



# Leaf Node Scan - Example 3



## Leaf Node Scan - Example 3



## 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 fail-fast **no-wait** mode.
- B+Tree implementation must cope with failed latch acquisitions.



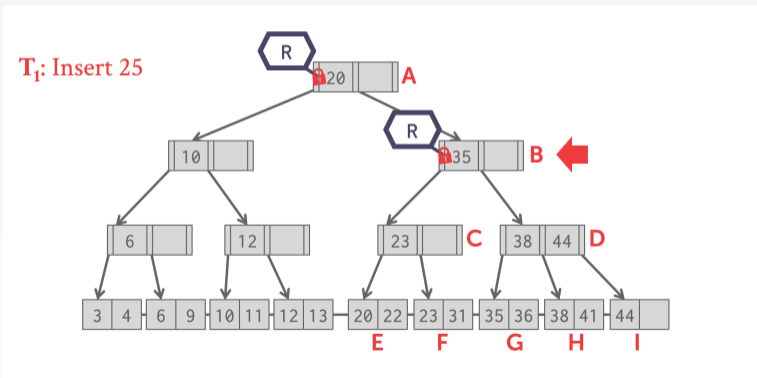
## *B*<sup>link</sup>-Tree

---

- Every time a leaf node overflows, we must update at least **three** nodes.
  - ▶ The leaf node being split.
  - ▶ The new leaf node being created.
  - ▶ The parent node.
- **Optimization:** When a leaf node overflows, delay updating its parent node.
- **Reference**

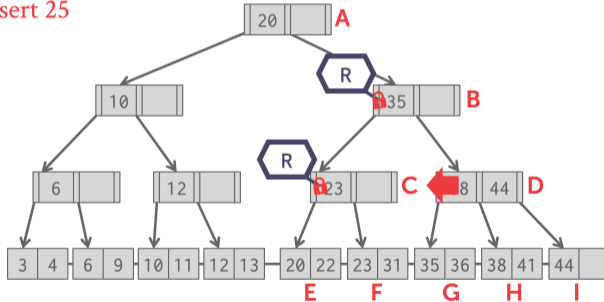


# B<sup>link</sup>-Tree Example

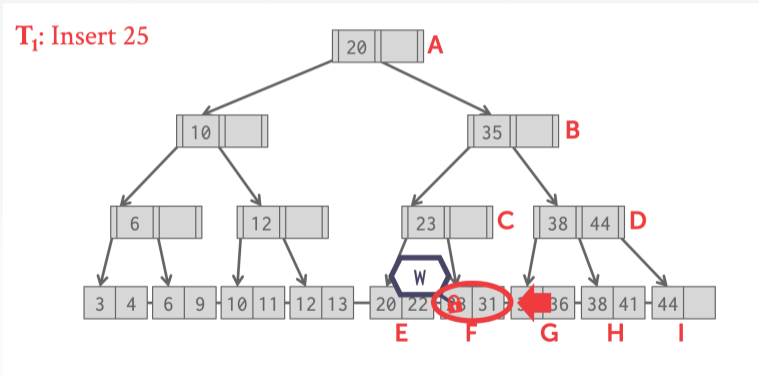


# $B^{link}$ -Tree Example

$T_1$ : Insert 25

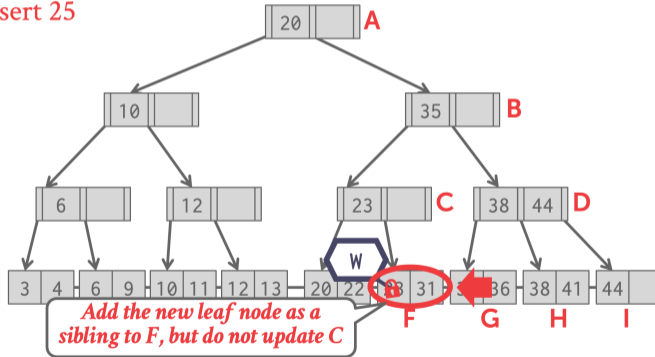


# $B^{link}$ -Tree Example



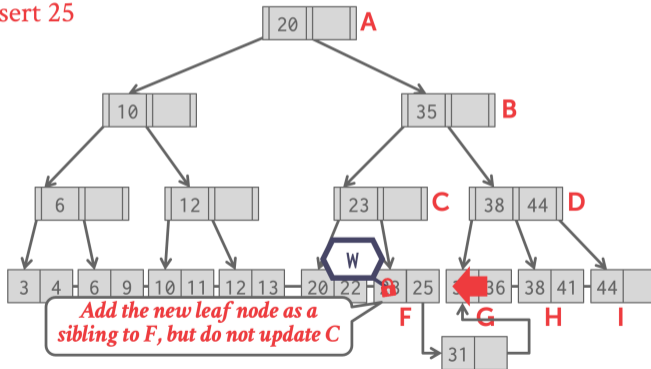
# Blink-Tree Example

$T_1$ : Insert 25

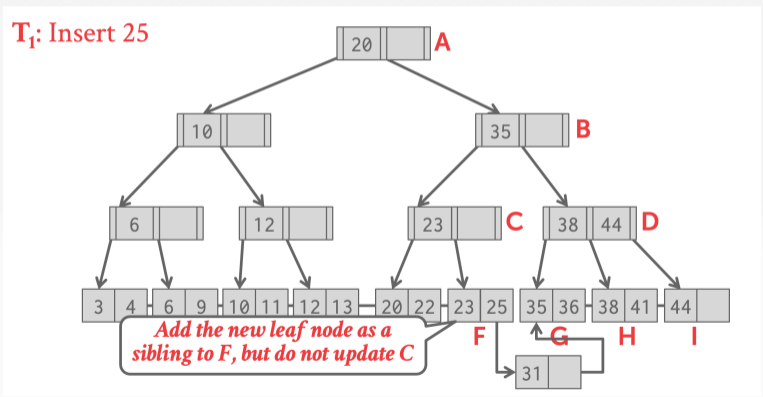


# B<sup>link</sup>-Tree Example

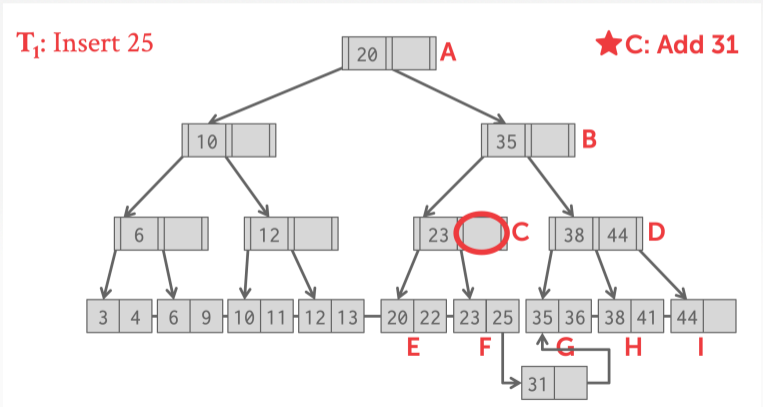
T<sub>1</sub>: Insert 25



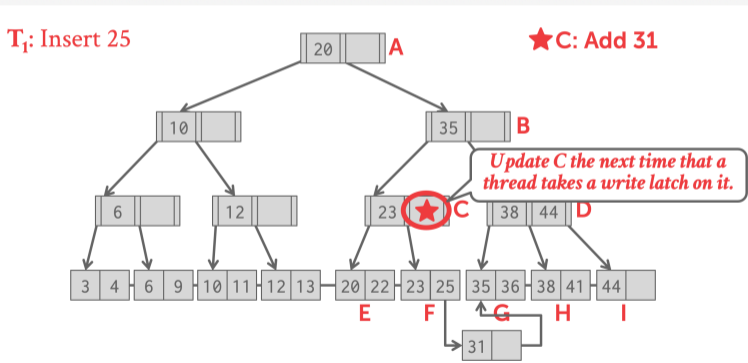
# B<sup>link</sup>-Tree Example



# Blink-Tree Example

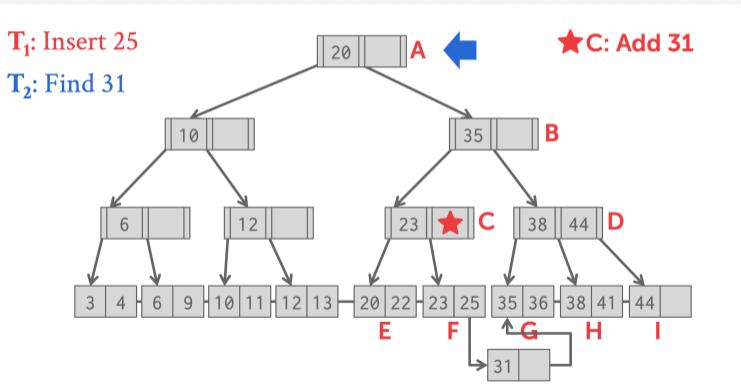


# Blink-Tree Example

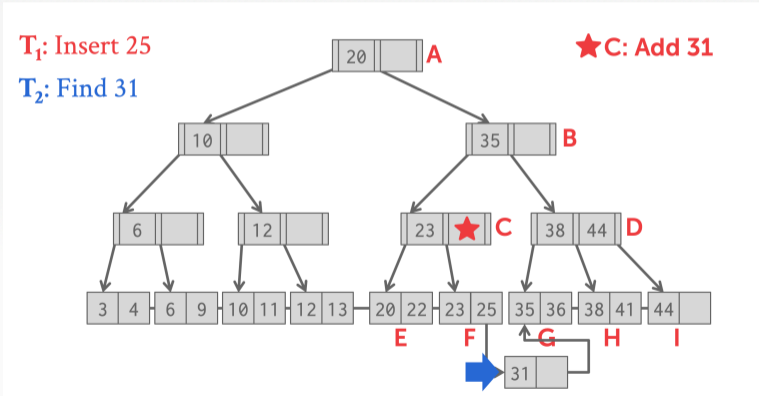




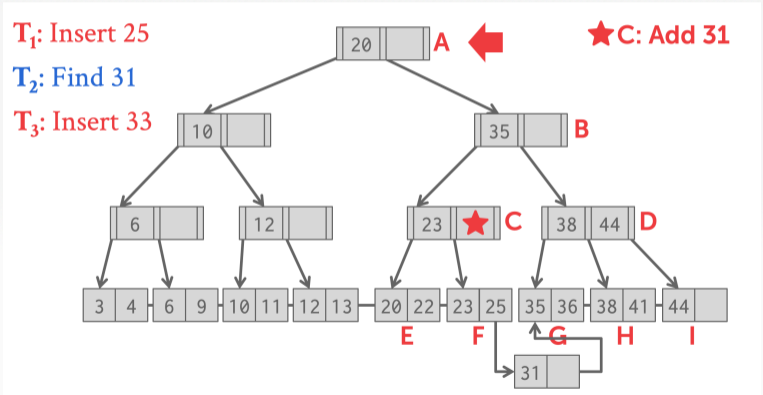
# $B^{\text{link}}$ -Tree Example



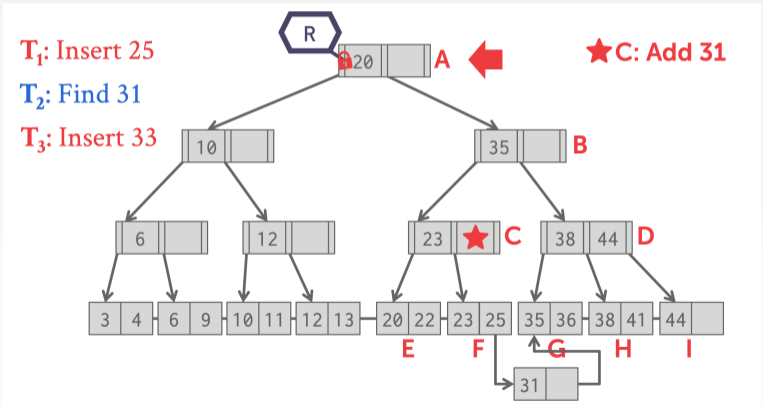
# Blink-Tree Example



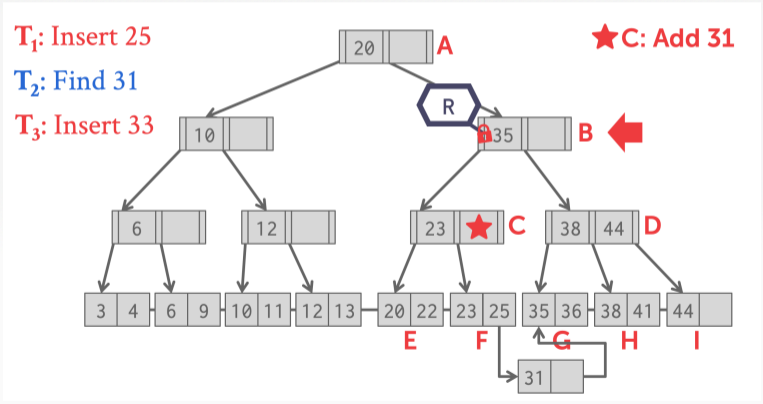
# B<sup>link</sup>-Tree Example



# Blink-Tree Example

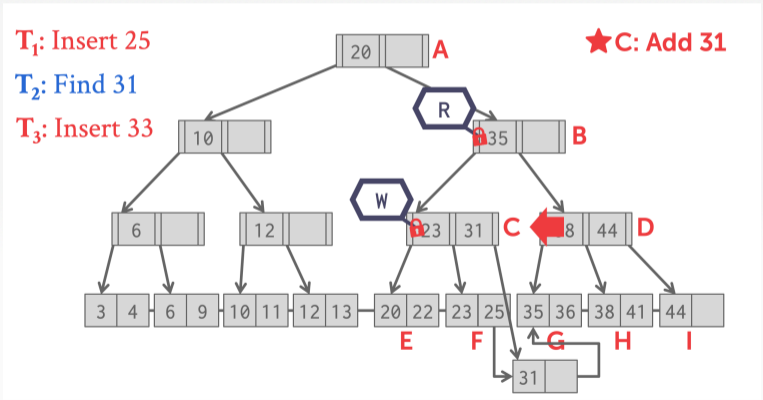


# Blink-Tree Example





# Blink-Tree Example







# 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 will learn about modern access methods.