

Index Concurrency Control

Administrivia

Open-book exam
No Homework

- Assignment 3 is due on Oct 19th @ 11:59pm
- Exercise Sheet 3 is due on Oct 19th @ 11:59pm (no late days allowed)

Sample Quiz on Carway

Recap

Index Data Structures

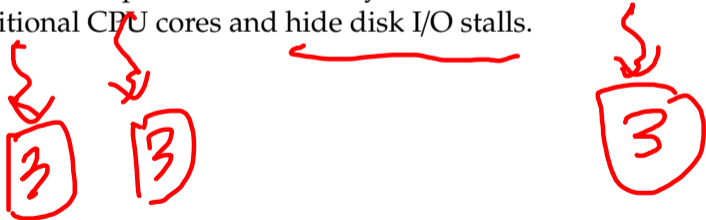
- List of Data Structures: Hash Tables, B+Trees, Radix Trees
- Most DBMSs automatically create an index to enforce integrity constraints.
- B+Trees are the way to go for indexing data.

/
SQL

Observation

B+ tree

- We assumed that all the data structures that we have discussed so far are single-threaded.
- But we need to allow multiple threads to safely access our data structures to take advantage of additional CPU cores and hide disk I/O stalls.



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:** Am I reading the data that I am supposed to read?
 - ▶ **Physical Correctness:** Is the internal representation of the object sound?

trans
locking
isolation
level



Today's Agenda

2nd lab

- Latches Overview
- Hash Table Latching
- B+Tree Latching
- Leaf Node Scans
- Delayed Parent Updates (B^{link}-Tree)

thread-safe

Latches Overview

Locks vs. Latches

• Locks

- ▶ Protects the database's logical contents from other txns.
- ▶ Held for the duration of the transaction.
- ▶ Need to be able to rollback changes.

• Latches

- ▶ Protects the critical sections of the DBMS's internal physical data structures from other threads.
- ▶ Held for the duration of the operation.
- ▶ Do not need to be able to rollback changes.

0 - SA: J₁
 0 - CA: r₂

multiple

~~locks~~ C++/OS jargon

Locks vs. Latches

	Locks	Latches
Separate...	User transactions	Threads
Protect...	Database Contents	In-Memory Data Structures
During...	Entire Transactions	Critical Sections
Modes...	Shared, Exclusive, Update, Intention	Read, Write (<i>a.k.a.</i> , Shared, Exclusive)
<u>Deadlock</u>	<u>Detection & Resolution</u>	<u>Avoidance</u>
...by...	Waits-for, Timeout, Aborts	Coding Discipline
Kept in...	Lock Manager	Protected Data Structure

Reference

(RW | SX)

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 holds the latch in any mode.

	Read	Write
Read	✓	X
Write	X	X

Compatibility Matrix

Latch Implementations

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

Latch Implementations

FUTEX

- Approach 1: Blocking OS Mutex

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

```
std::mutex m;
```

```
m.lock();
```

```
// Do something special...
```

```
m.unlock();
```

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?  
}
```

lock structure

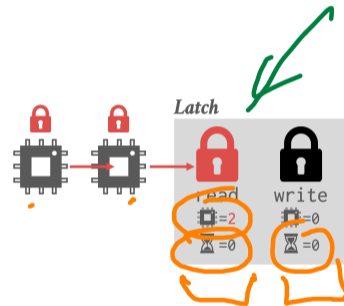
Compare &
Swap
(CAS)

0 → 1 (set)
1 → 0 (clear)

Latch Implementations

- Approach 3: Reader-Writer Latch

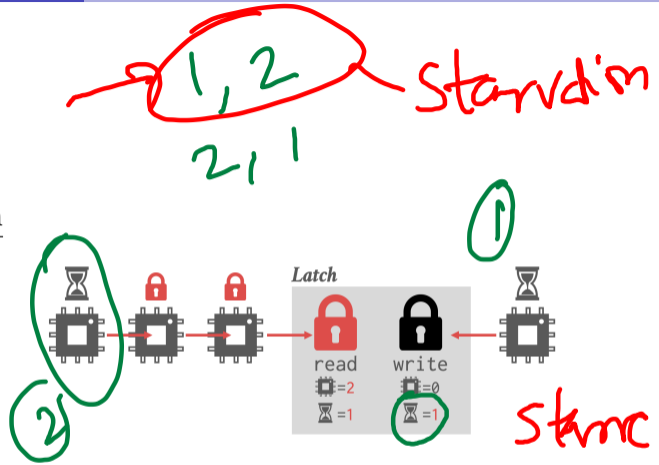
- ▶ Allows for concurrent readers
- ▶ Must manage read/write queues to avoid starvation
- ▶ Can be implemented on top of spinlocks



Latch Implementations

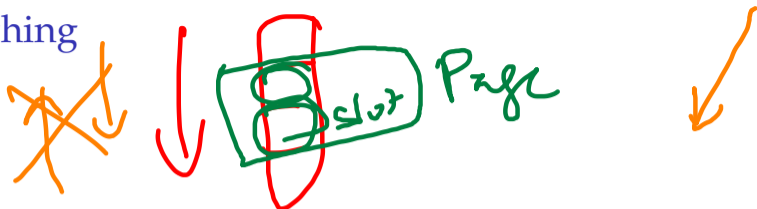
• Approach 3: Reader-Writer Latch

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



Hash Table Latching

Hash Table Latching



- Easy to support concurrent access due to the limited ways in which 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 latch on the entire table (*i.e.*, 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.

R, W

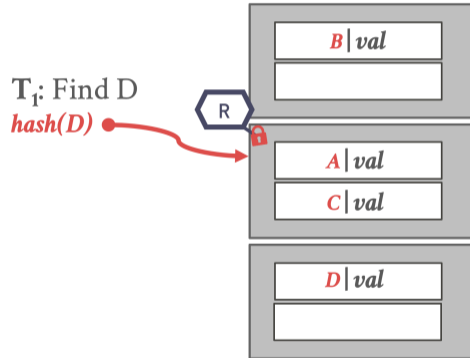
- Approach 2: Slot Latches

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

O/I

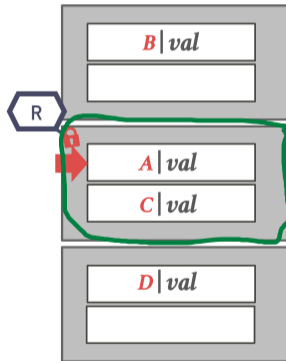
more consistent

Hash Table - Page Latches



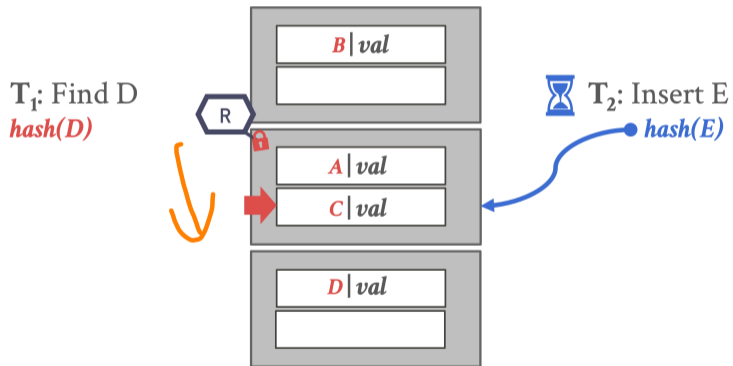
Hash Table - Page Latches

T_1 : Find D
 $hash(D)$

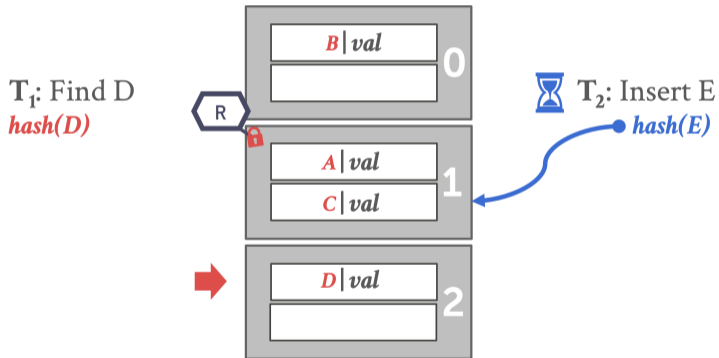


T_2 : Insert E
 $hash(E)$

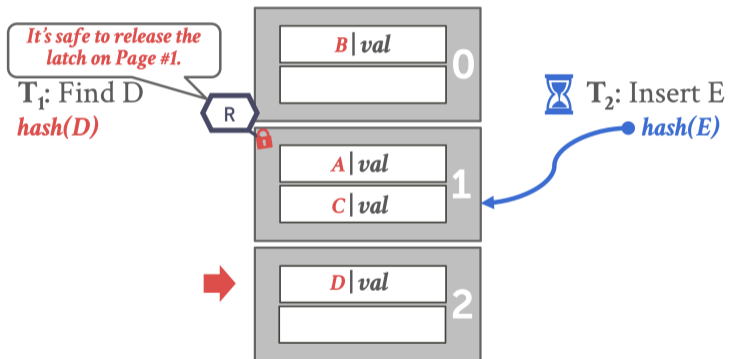
Hash Table - Page Latches



Hash Table - Page Latches

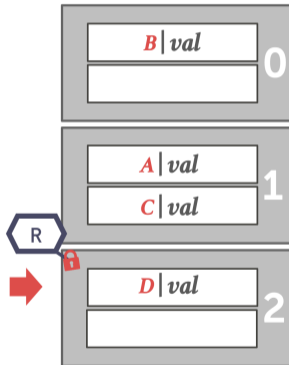


Hash Table - Page Latches



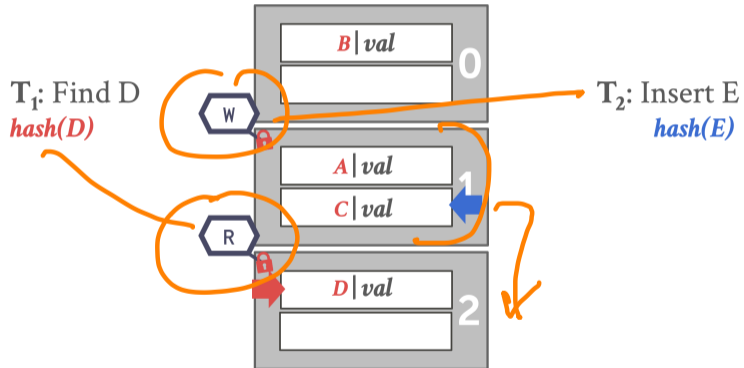
Hash Table - Page Latches

T_1 : Find D
 $hash(D)$

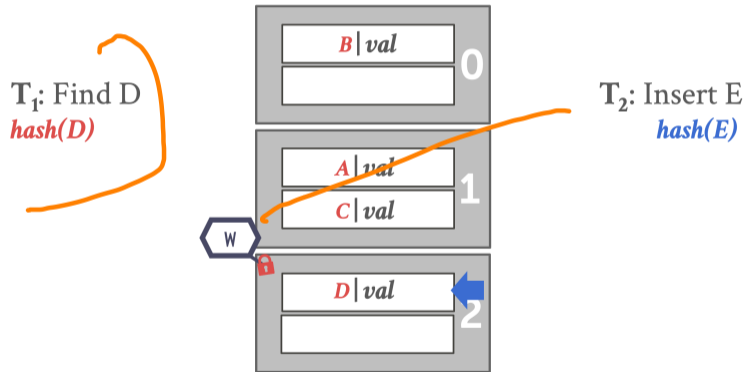


T_2 : Insert E
 $hash(E)$

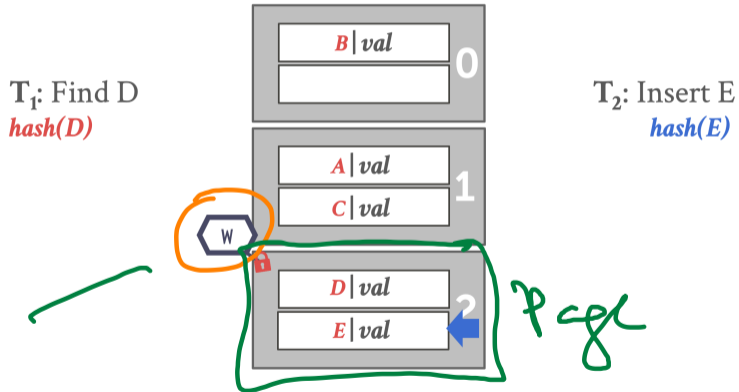
Hash Table - Page Latches



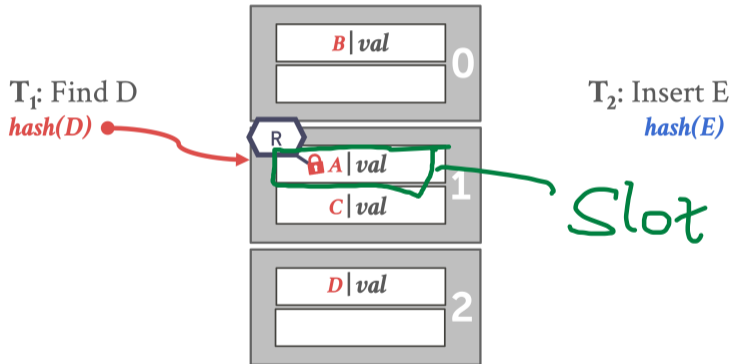
Hash Table - Page Latches



Hash Table - Page Latches

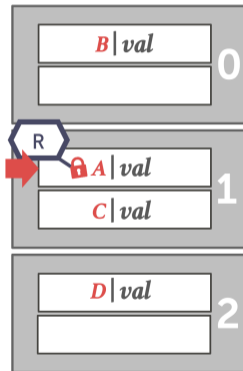


Hash Table - Slot Latches



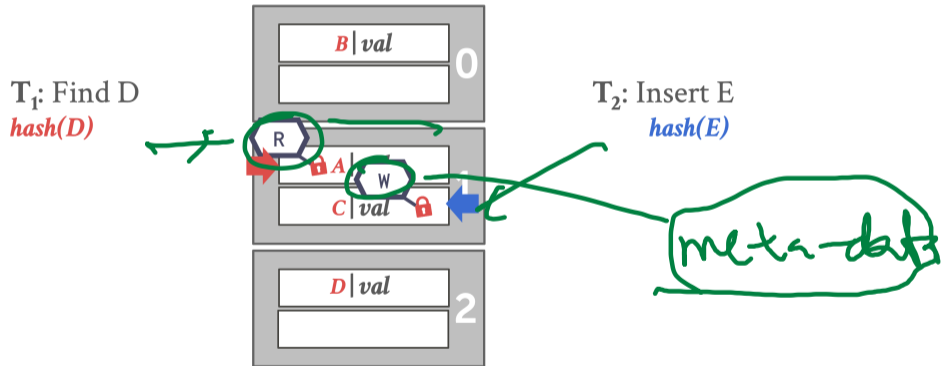
Hash Table - Slot Latches

T_1 : Find D
 $hash(D)$

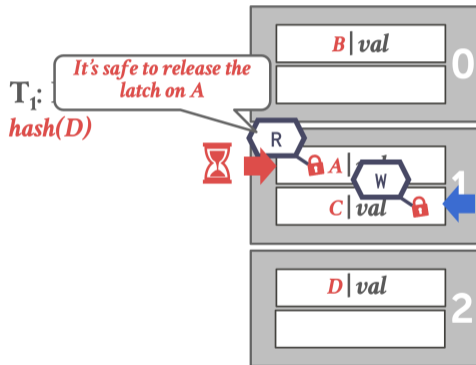


T_2 : Insert E
 $hash(E)$

Hash Table - Slot Latches



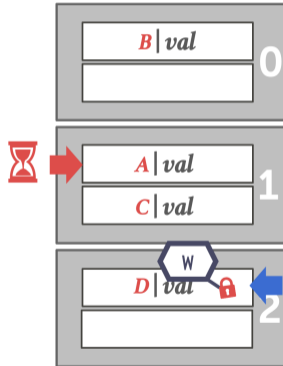
Hash Table - Slot Latches



T_2 : Insert E
 $hash(E)$

Hash Table - Slot Latches

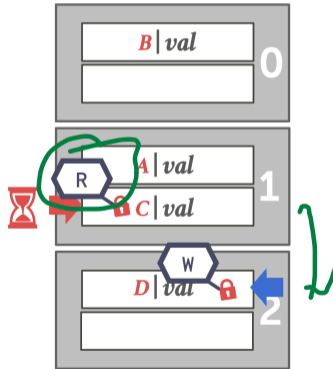
T_1 : Find D
hash(D)



T_2 : Insert E
hash(E)

Hash Table - Slot Latches

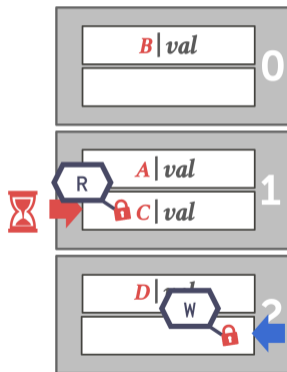
T_1 : Find D
hash(D)



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

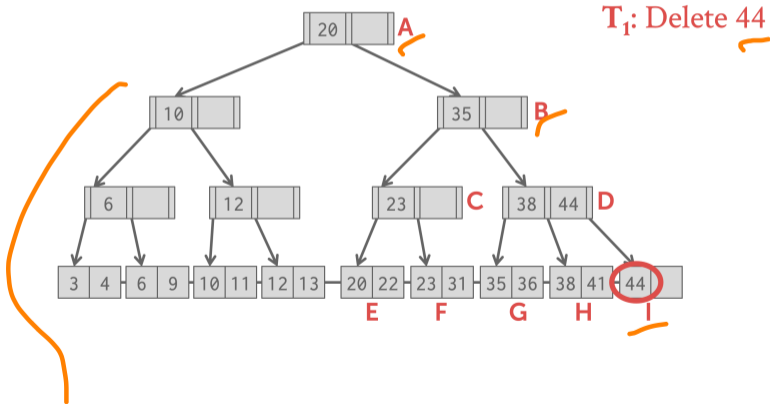
B+Tree Concurrency Control

- 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.

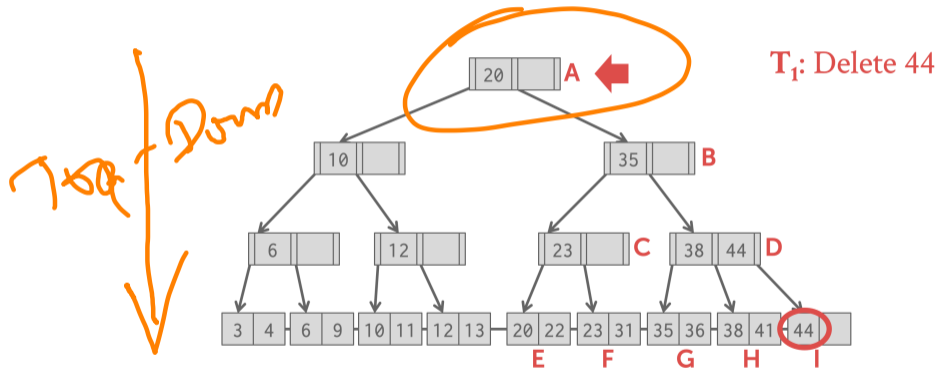
Single
node

→ (traversing)

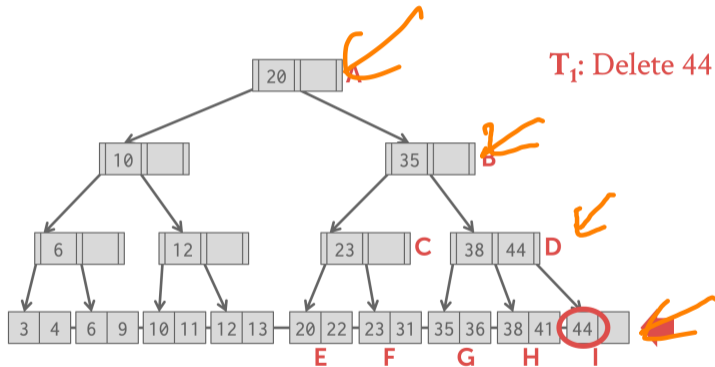
B+Tree Concurrency Control: Example



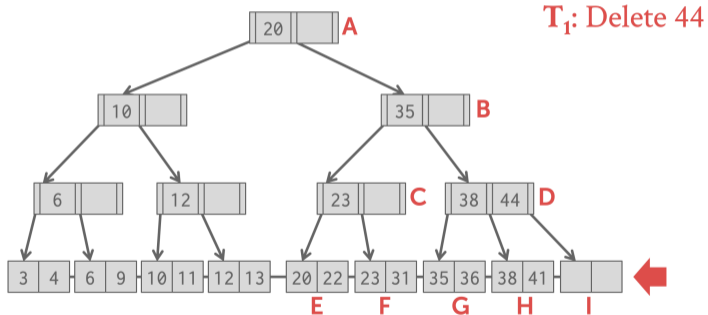
B+Tree Concurrency Control: Example



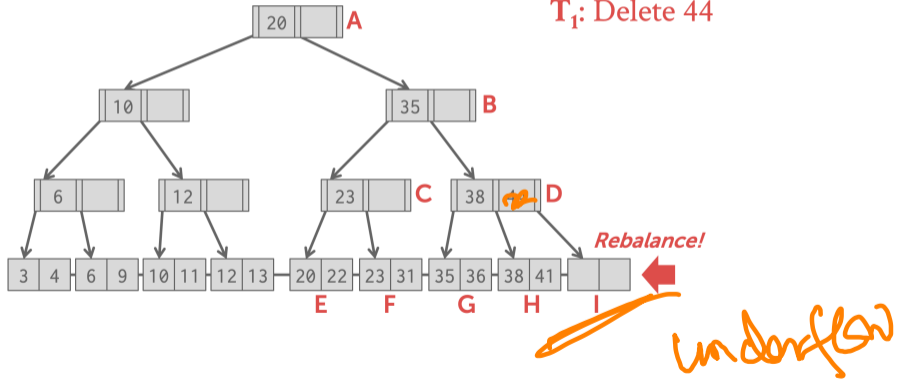
B+Tree Concurrency Control: Example



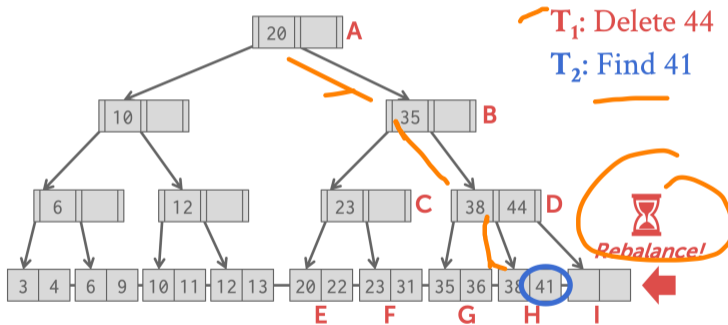
B+Tree Concurrency Control: Example



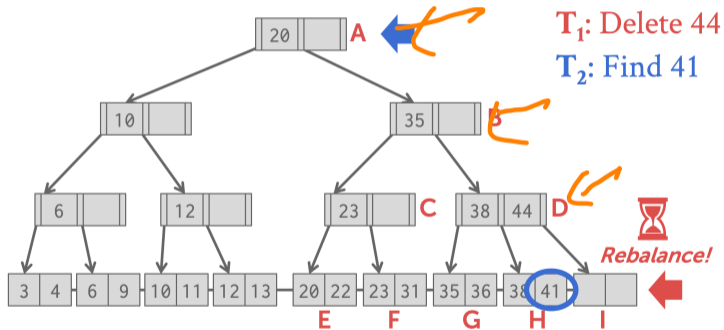
B+Tree Concurrency Control: Example



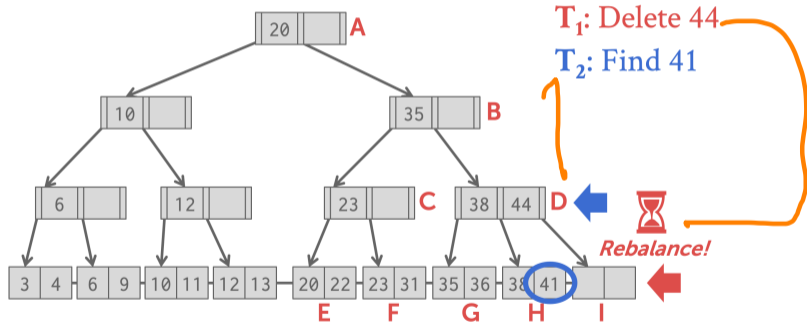
B+Tree Concurrency Control: Example



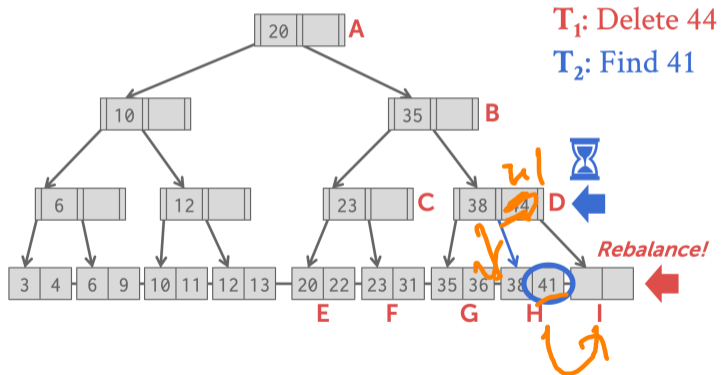
B+Tree Concurrency Control: Example



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)

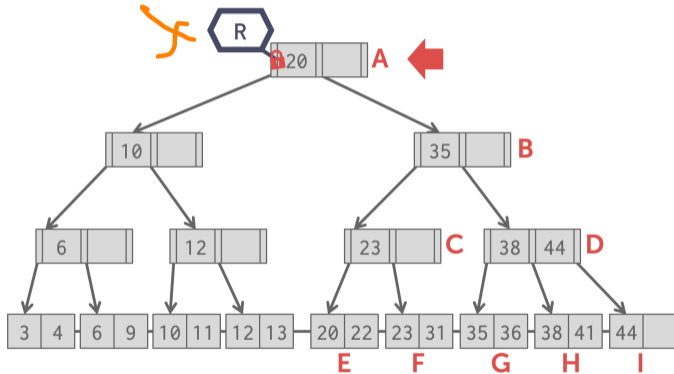
child

Latch Crabbing/Coupling

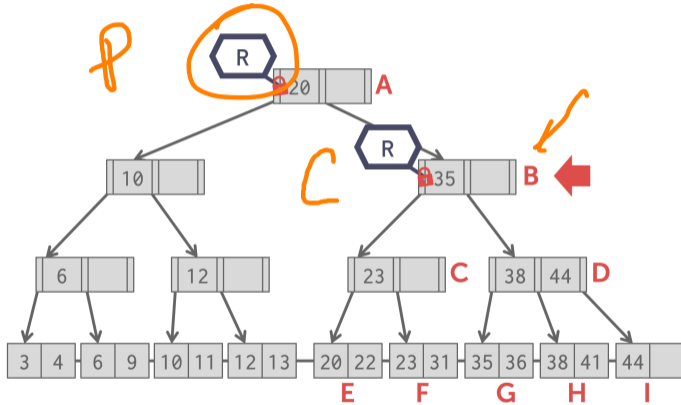
- **Find:** Start at root and go down; repeatedly,
 - ▶ Acquire **R** latch on child
 - ▶ Then unlatch parent
- **Insert/Delete:** Start at root and go down, obtaining **W** latches as needed. Once child is latched, check if it is safe:
 - ▶ If child is safe, release all latches on ancestors.

Examples

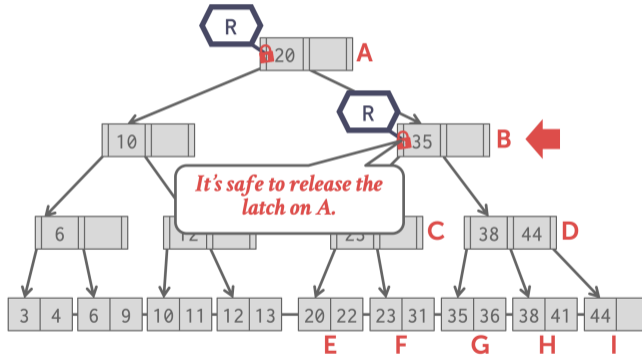
Example 1 - Find 38



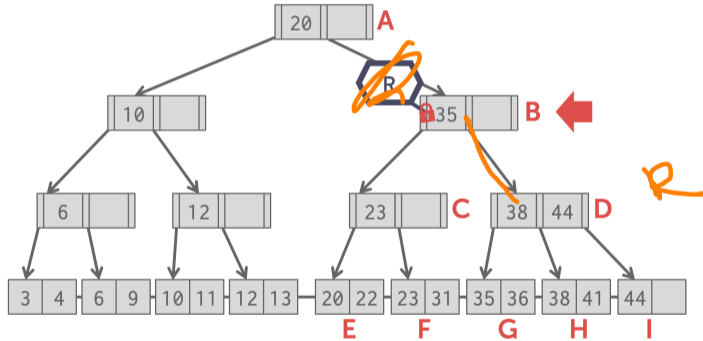
Example 1 - Find 38



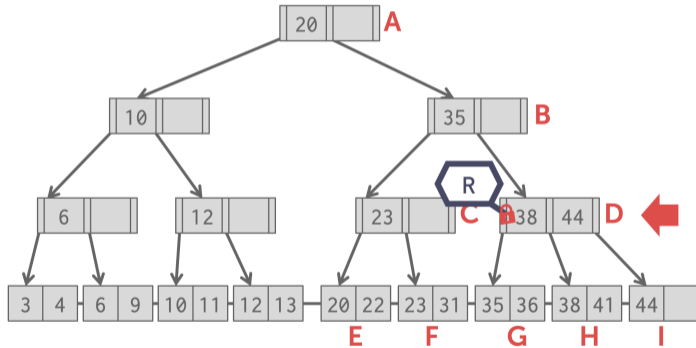
Example 1 - Find 38



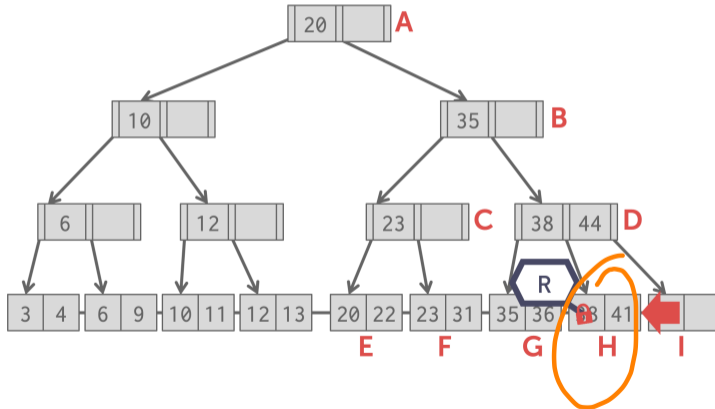
Example 1 - Find 38



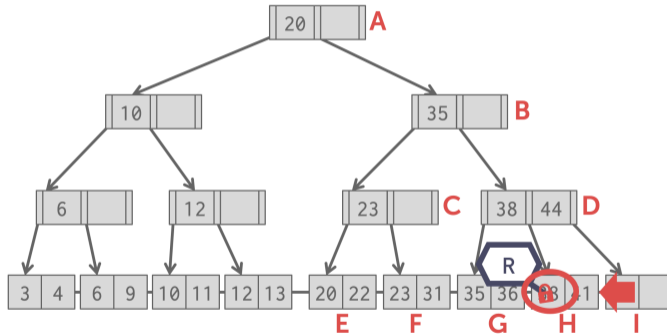
Example 1 - Find 38



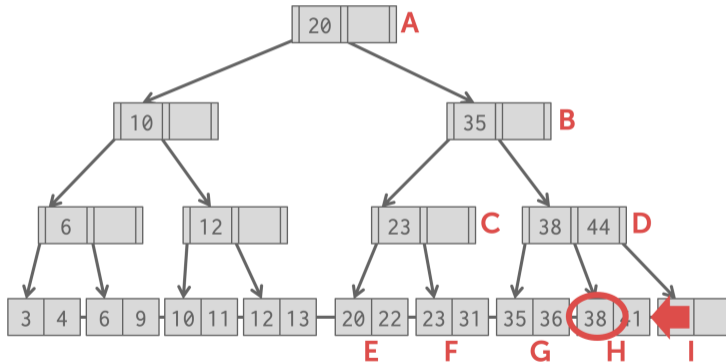
Example 1 - Find 38



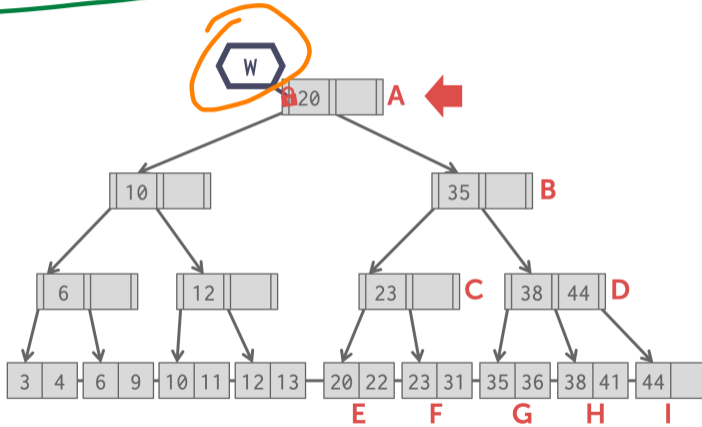
Example 1 - Find 38



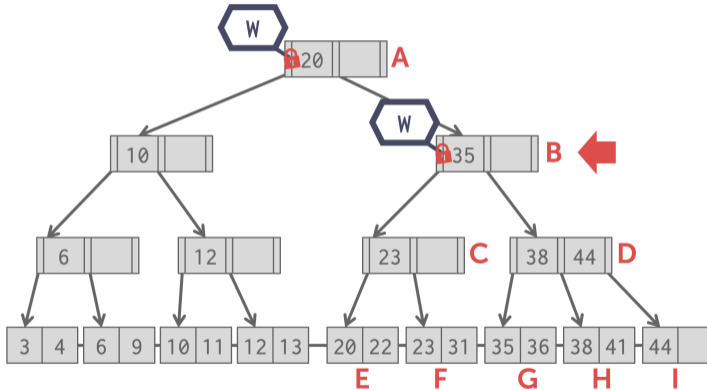
Example 1 - Find 38



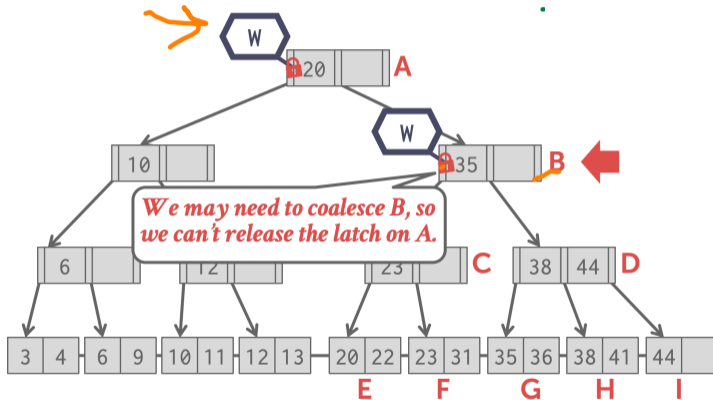
Example 2 - Delete 38



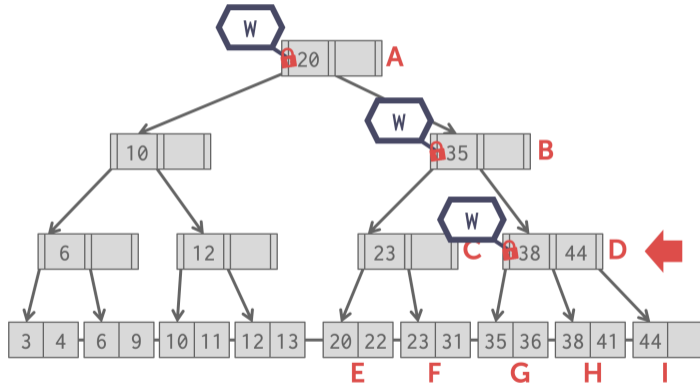
Example 2 - Delete 38



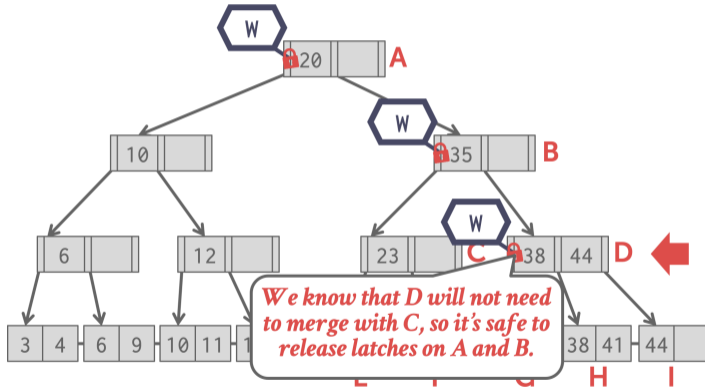
Example 2 - Delete 38



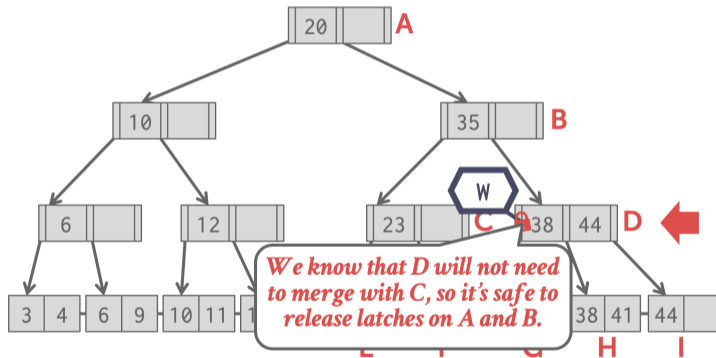
Example 2 - Delete 38



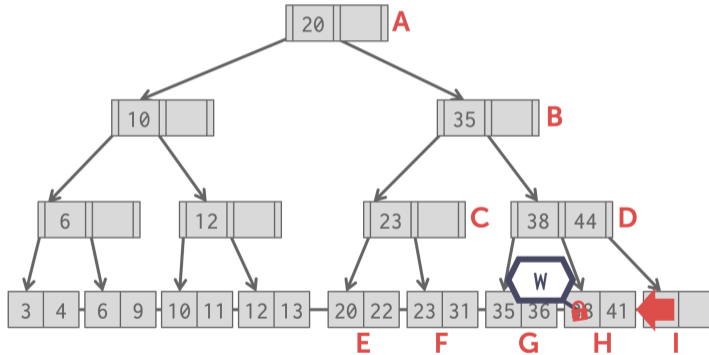
Example 2 - Delete 38



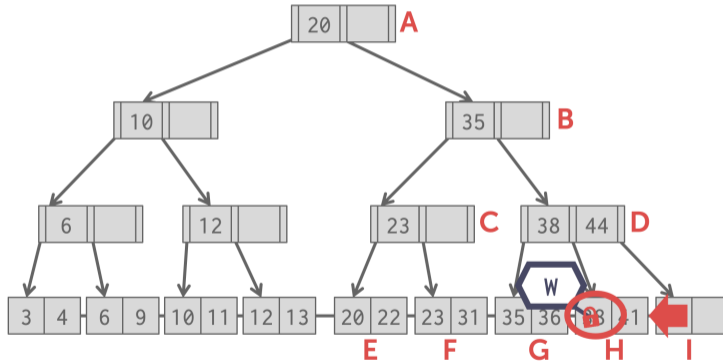
Example 2 - Delete 38



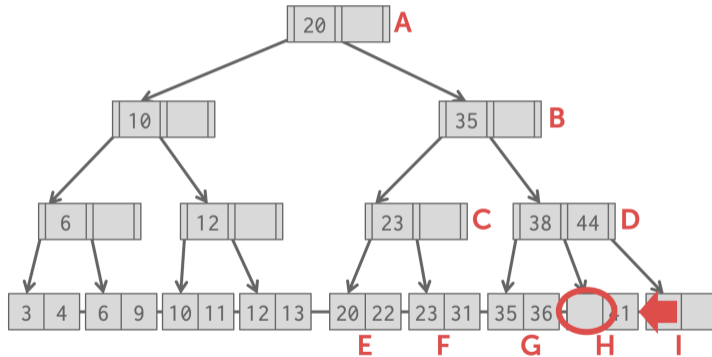
Example 2 - Delete 38



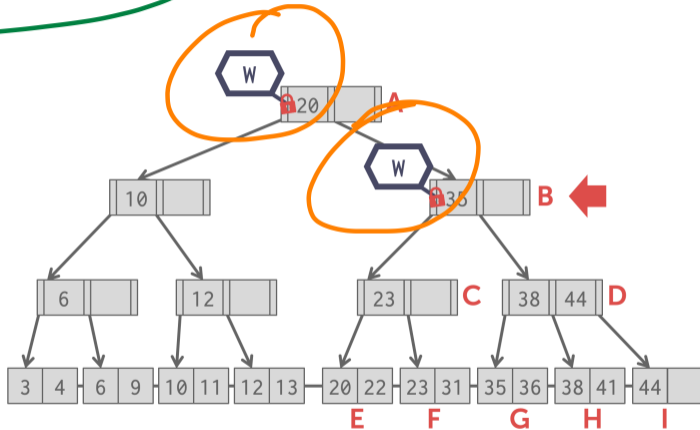
Example 2 - Delete 38



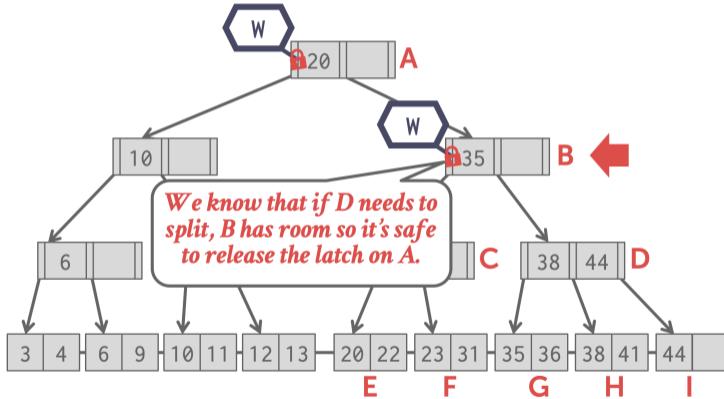
Example 2 - Delete 38



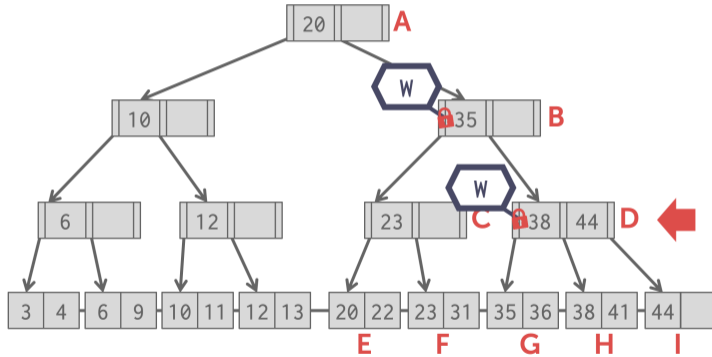
Example 3 - Insert 45



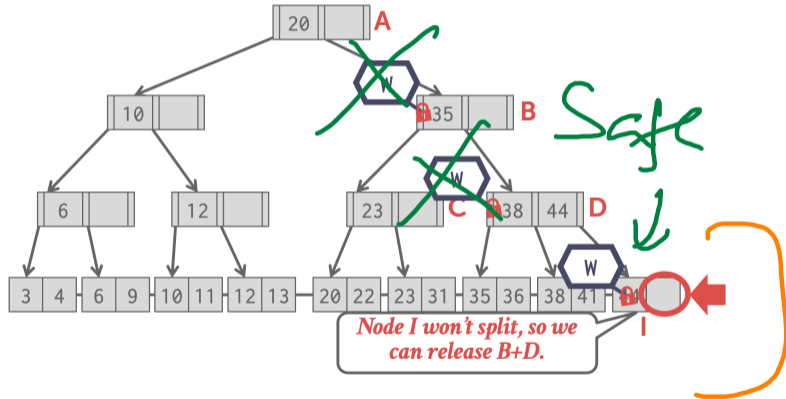
Example 3 - Insert 45



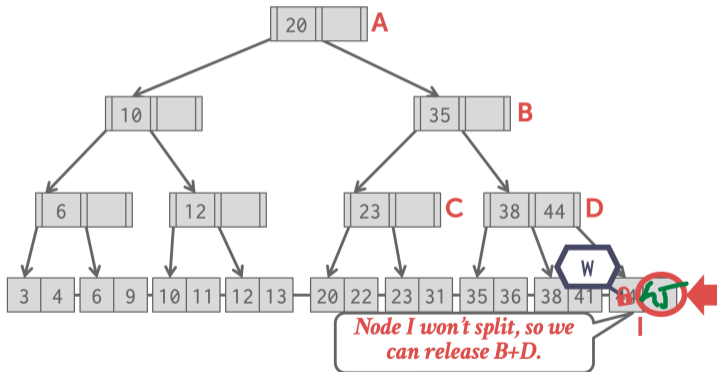
Example 3 - Insert 45



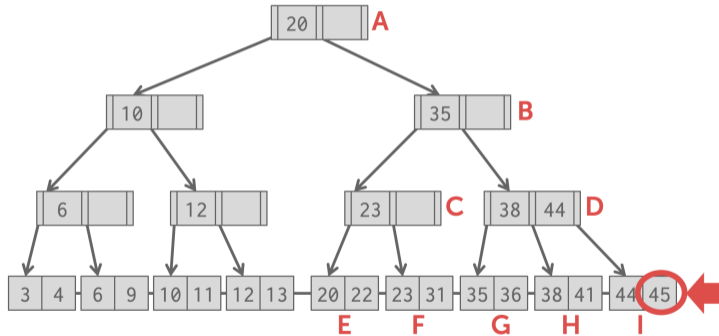
Example 3 - Insert 45



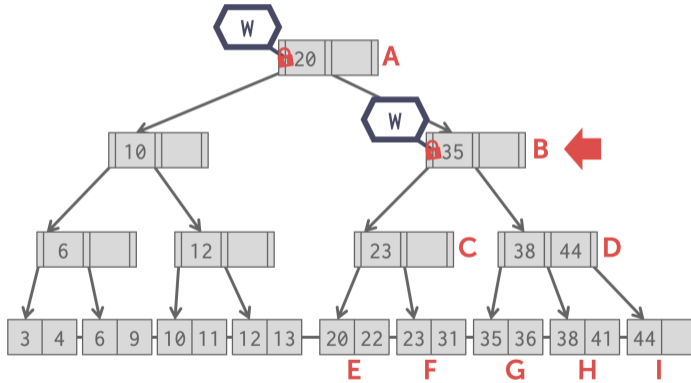
Example 3 - Insert 45



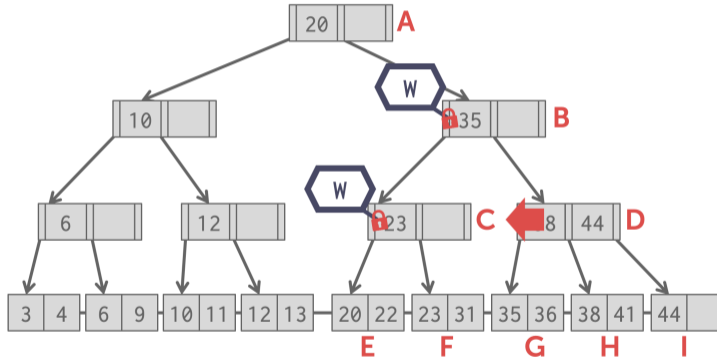
Example 3 - Insert 45



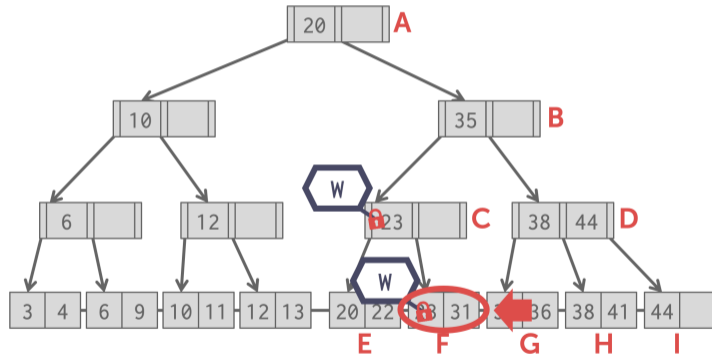
Example 4 - Insert 25



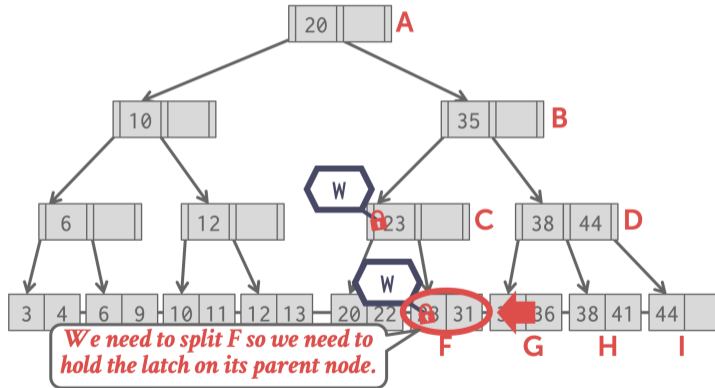
Example 4 - Insert 25



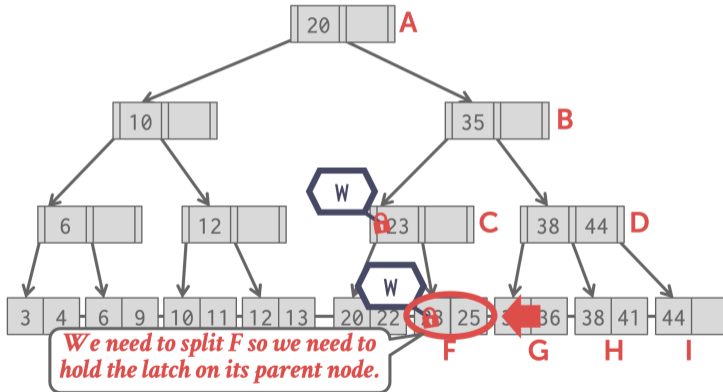
Example 4 - Insert 25



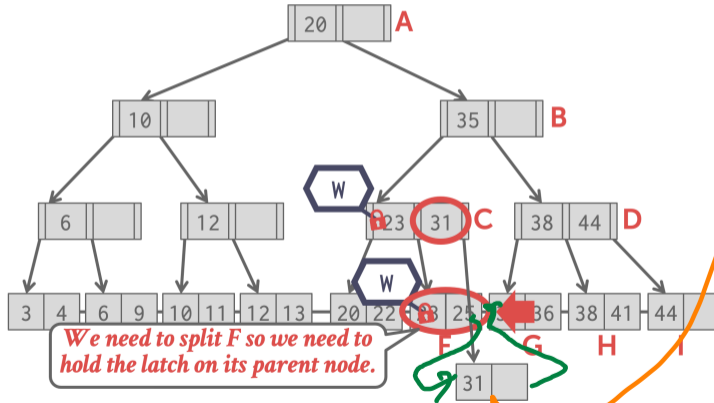
Example 4 - Insert 25



Example 4 - Insert 25



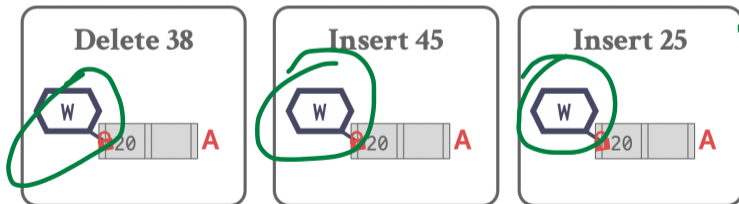
Example 4 - Insert 25



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?

Optimism vs Pessimism



Better Latching Algorithm

- Assume that the leaf node is safe.
- Use read latches and crabbing to reach it, and then verify that it is safe.
- If leaf is not safe, then do previous algorithm using write latches.

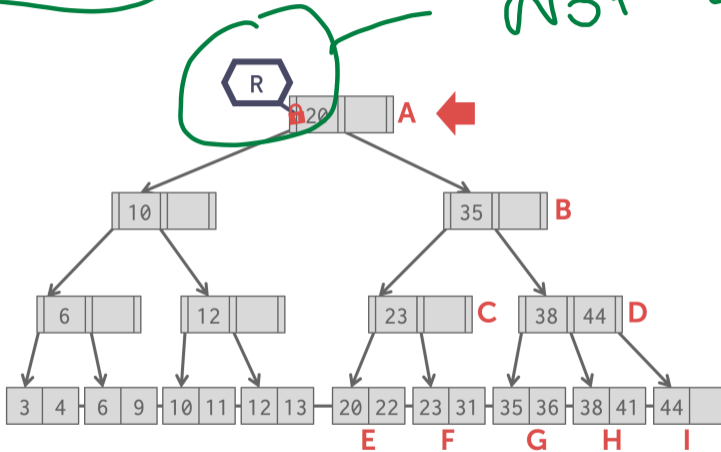
• Reference

contention is low
in the common case

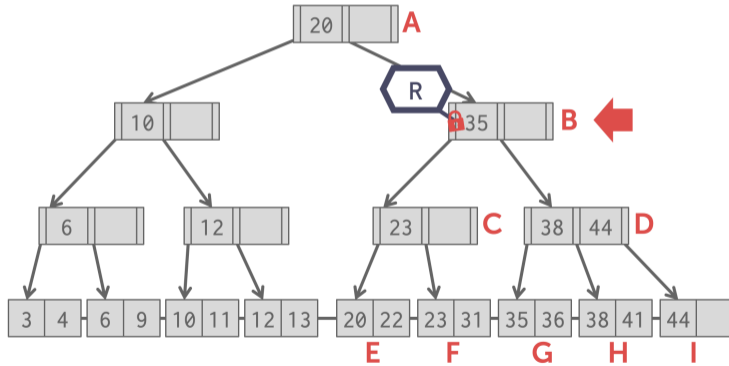
fallback to slower algorithm if needed

Example 2 - Delete 38

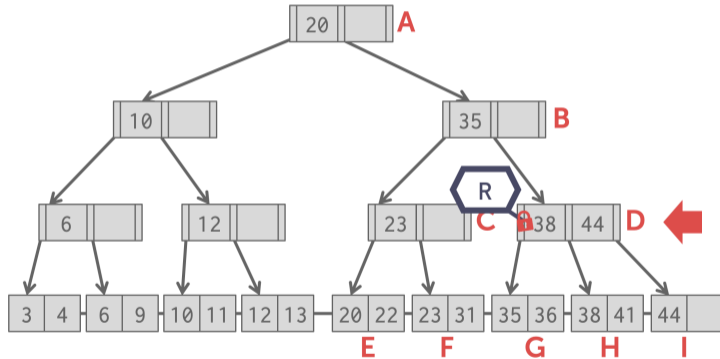
Not W



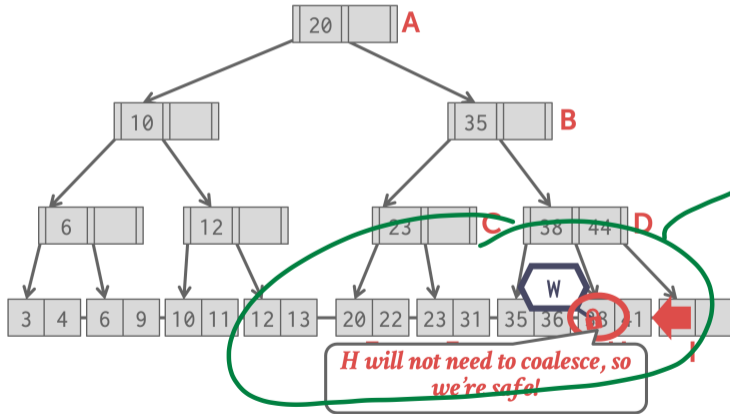
Example 2 - Delete 38



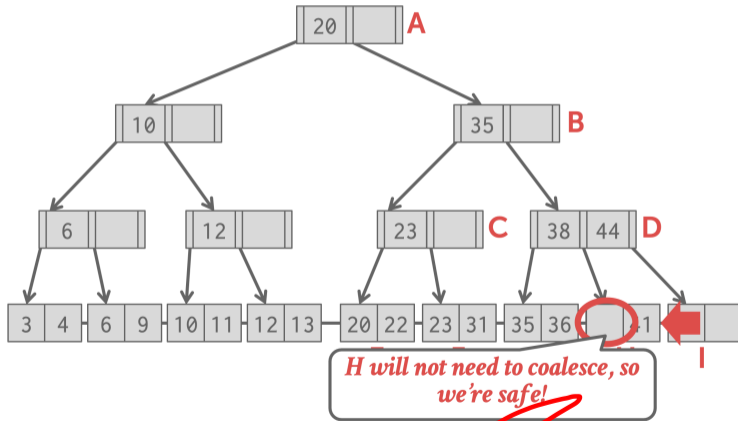
Example 2 - Delete 38



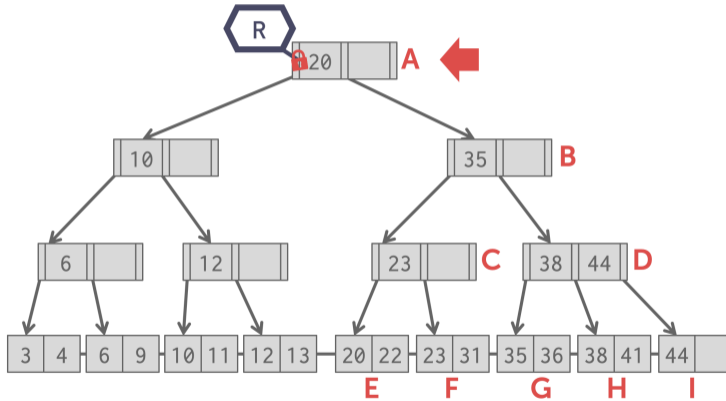
Example 2 - Delete 38



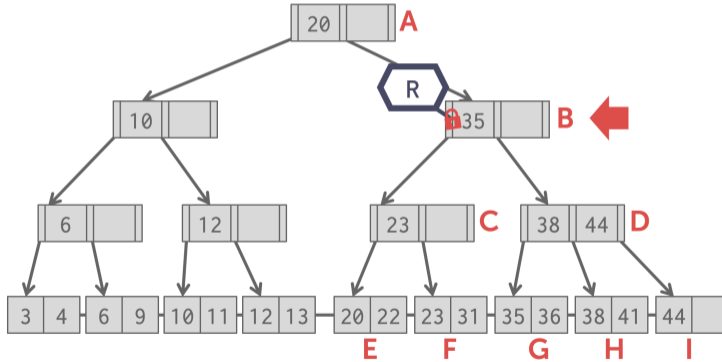
Example 2 - ~~Insert 25~~



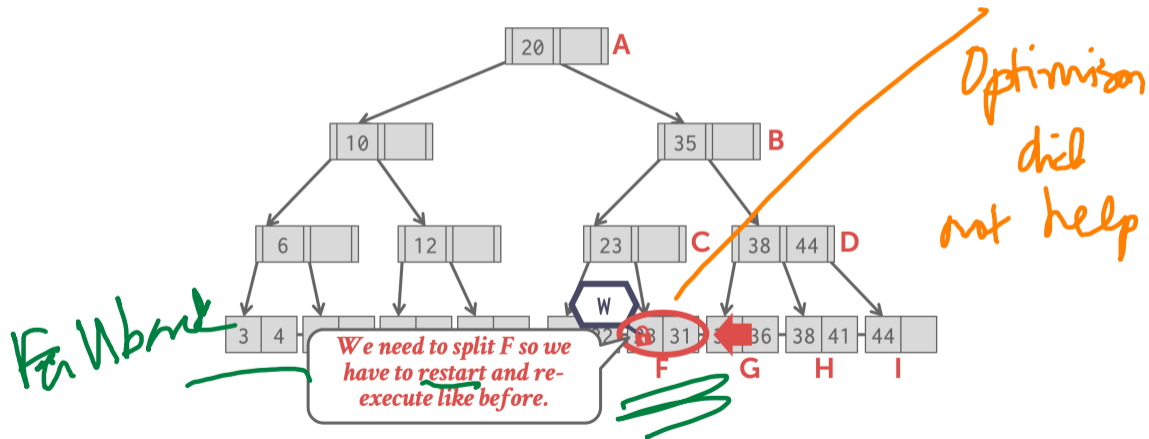
Example 4 - Insert 25



Example 4 - Insert 25



Example 4 - Insert 25



Better Latching Algorithm

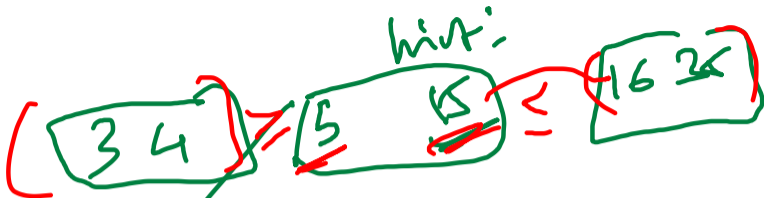
Optimistic

- **Find:** 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 **W** 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.

not the common case

Leaf Node Scans

Observation



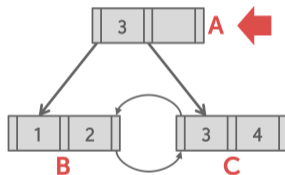
- The threads in all the examples so far have acquired latches in a top-down manner.
 - ▶ A thread can only acquire a latch from a node that is below its current node.
 - ▶ If the desired latch is unavailable, the thread must wait until it becomes available.
- But what if we want to move from one leaf node to another leaf node?
- Leaf nodes can include hint keys to approximate the next key at your sibling.

Range Scan

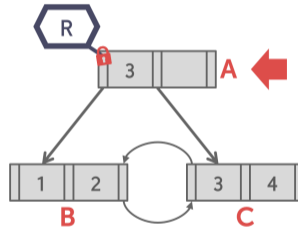
Salary > 100k < 200

Leaf Node Scan - Example 1

T_1 : Find Keys < 4



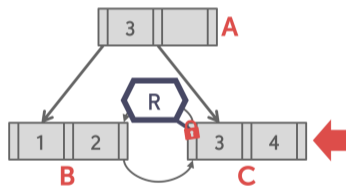
Leaf Node Scan - Example 1



T_1 : Find Keys < 4

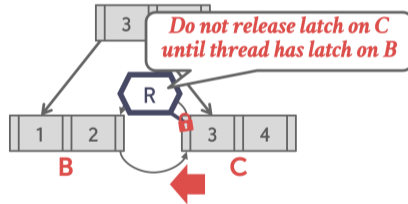
Leaf Node Scan - Example 1

T_1 : Find Keys < 4



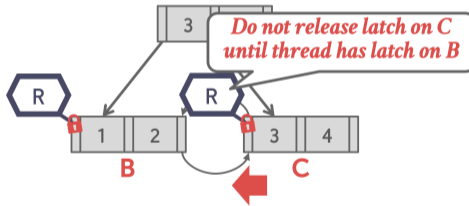
Leaf Node Scan - Example 1

T_1 : Find Keys < 4



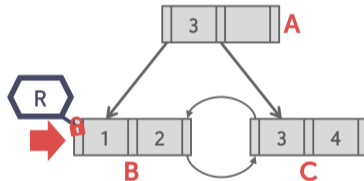
Leaf Node Scan - Example 1

T_1 : Find Keys < 4

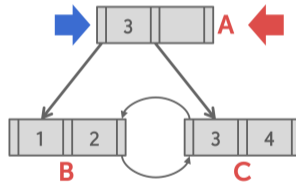


Leaf Node Scan - Example 1

T_1 : Find Keys < 4



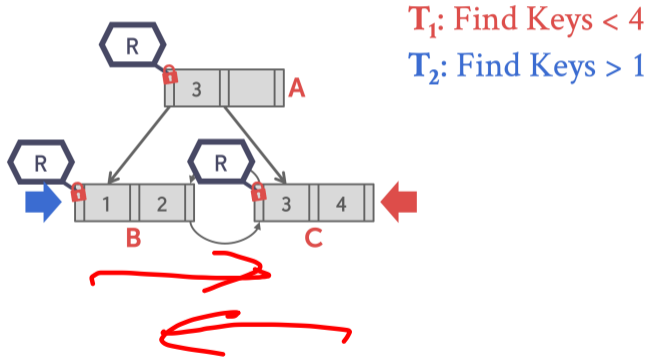
Leaf Node Scan - Example 2



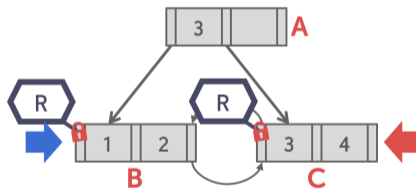
T_1 : Find Keys < 4

T_2 : Find Keys > 1

Leaf Node Scan - Example 2



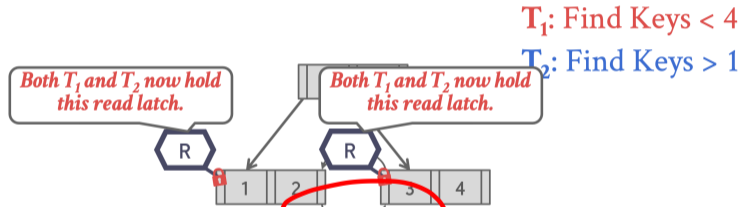
Leaf Node Scan - Example 2



T_1 : Find Keys < 4

T_2 : Find Keys > 1

Leaf Node Scan - Example 2

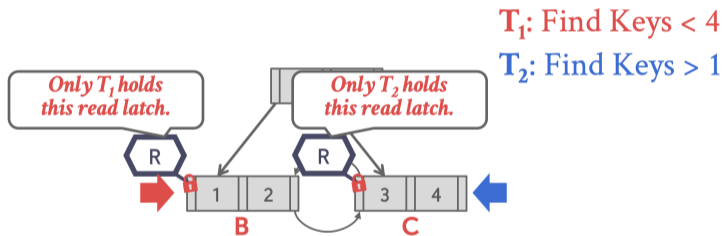


Discipline out met

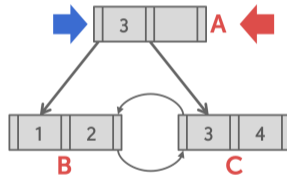
latching

X →

Leaf Node Scan - Example 2



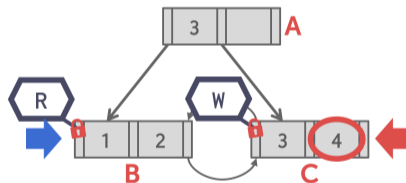
Leaf Node Scan - Example 3



T_1 : Delete 4

T_2 : Find Keys > 1

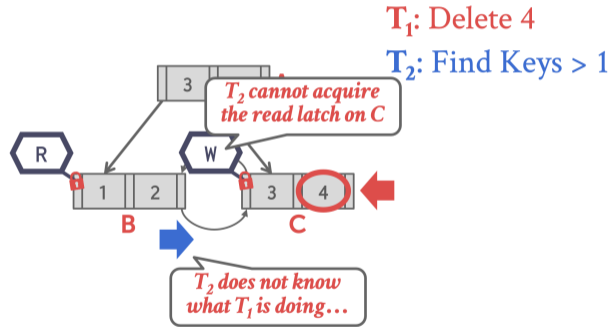
Leaf Node Scan - Example 3



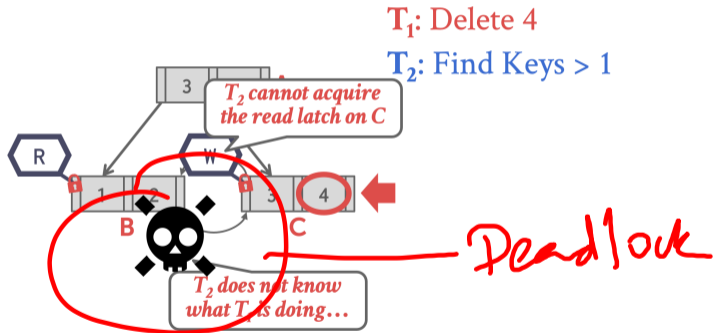
T_1 : Delete 4

T_2 : Find Keys > 1

Leaf Node Scan - Example 3



Leaf Node Scan - Example 3



Leaf Node Scans

Latches do

- 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.

Deadlock Avoidance

Blink-Tree

Delayed Parent Updates

Blink-Tree

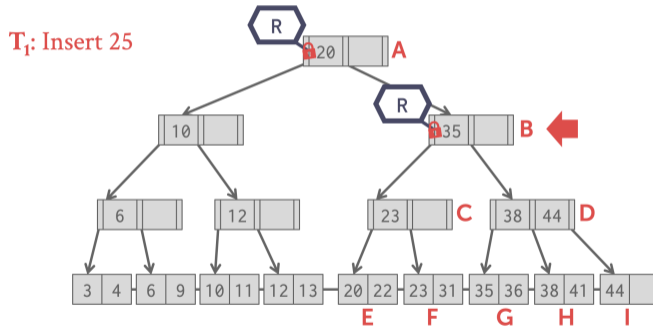


- Every time a leaf node overflows, we must update at least **three** nodes.
 - ▶ The leaf node being split.
 - ▶ The new leaf node being created.
 - ▶ The parent node.

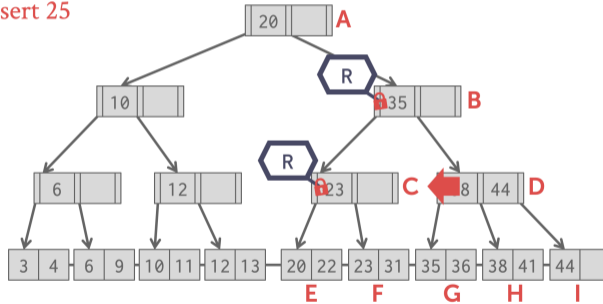
- Optimization: When a leaf node overflows, delay updating its parent node.

- Reference

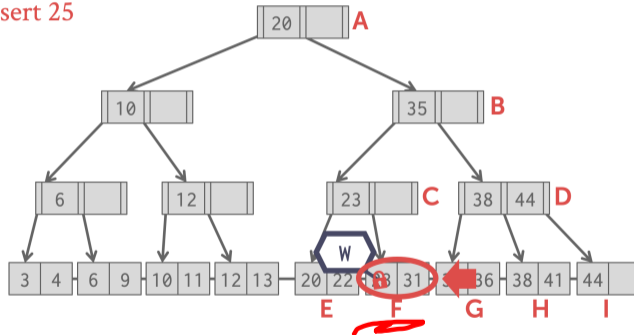
Blink-Tree Example



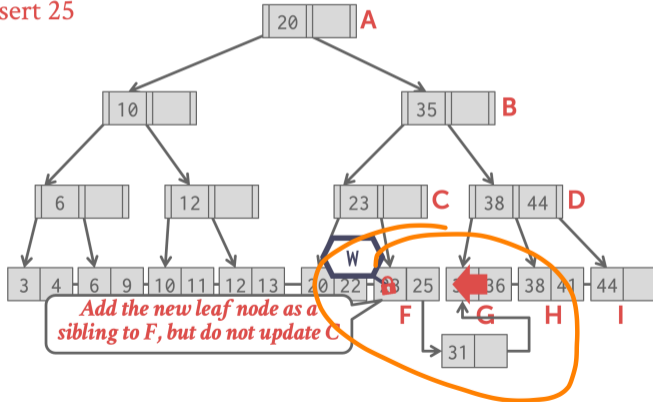
Blink-Tree Example

 T_1 : Insert 25

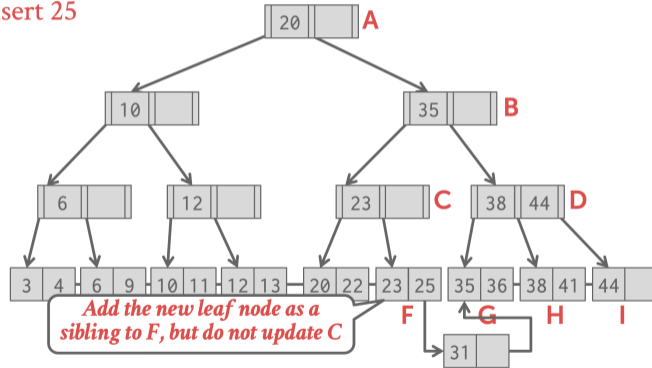
Blink-Tree Example

 T_i : Insert 25

Blink-Tree Example

 T_1 : Insert 25

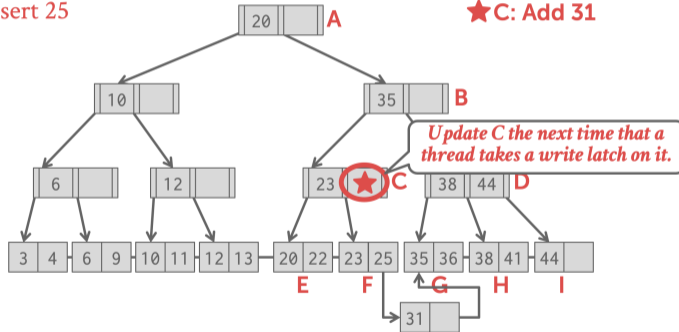
Blink-Tree Example

 T_1 : Insert 25

Blink-Tree Example

T₁: Insert 25

★C: Add 31

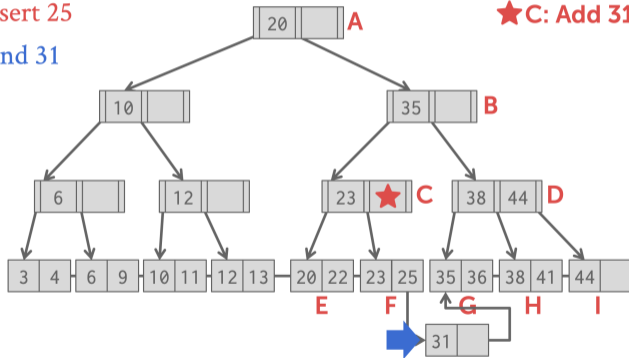


Blink-Tree Example

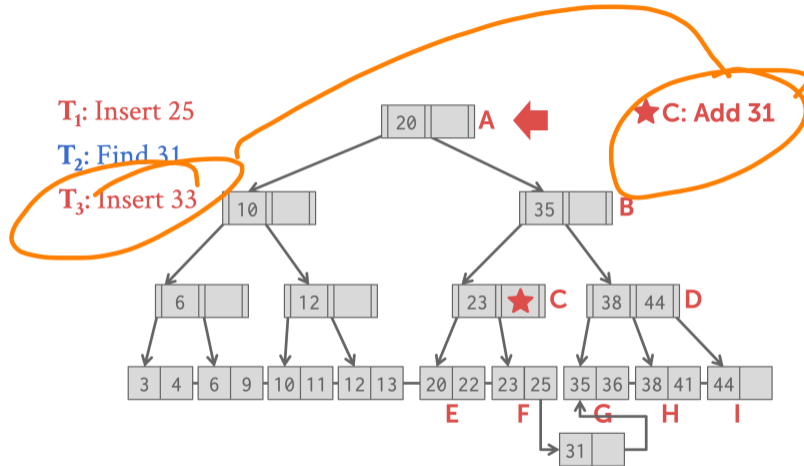
T_1 : Insert 25

T_2 : Find 31

★C: Add 31



Blink-Tree Example



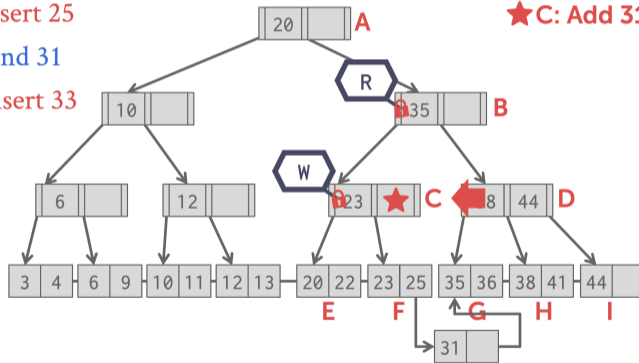
Blink-Tree Example

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33

★C: Add 31



Conclusion

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.

Buffer
manager lab