

CS 7270: Advanced Database Systems Fall 2025

Lecture 06

Concurrency control #1

Prashant Pandey

p.pandey@northeastern.edu

Acknowledgement: Slides taken from Prof. Andy Pavlo, CMU

OBSERVATION

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

LOCKS VS. LATCHES

Locks

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

Latches

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

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
Deadlock	Detection & Resolution	Avoidance
...by...	Waits-for, Timeout, Aborts	Coding Discipline
Kept in...	Lock Manager	Protected Data Structure

LOCKS VS. LATCHES

Next

Locks

Separate... User 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

Threads
In-Memory Data Structures
Critical Sections
Read, Write

Avoidance
Coding Discipline
Protected Data Structure

B+TREE CONCURRENCY CONTROL

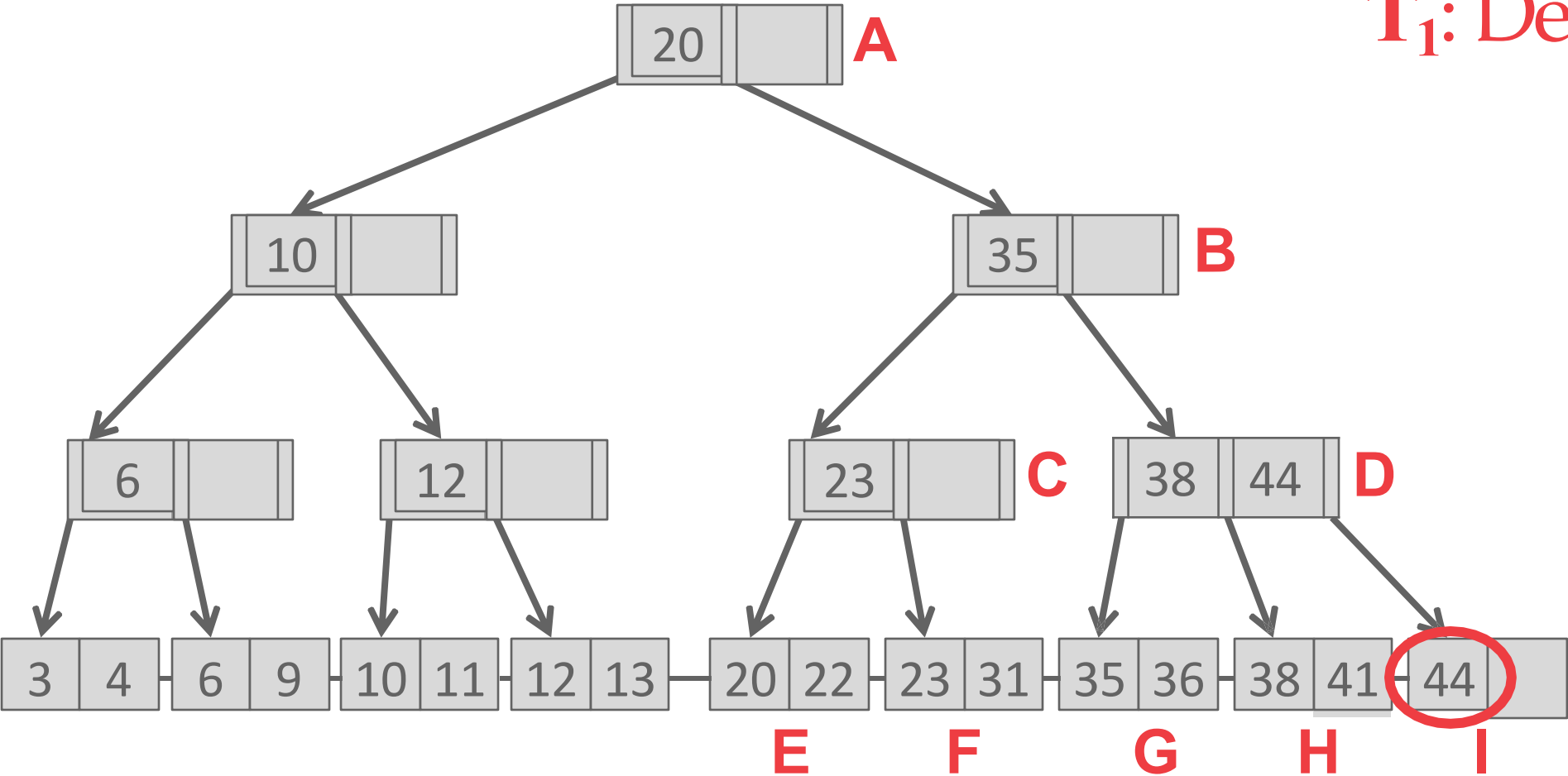
We want to allow multiple threads to read and update a B+Tree at the same time.

We need to protect from 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.

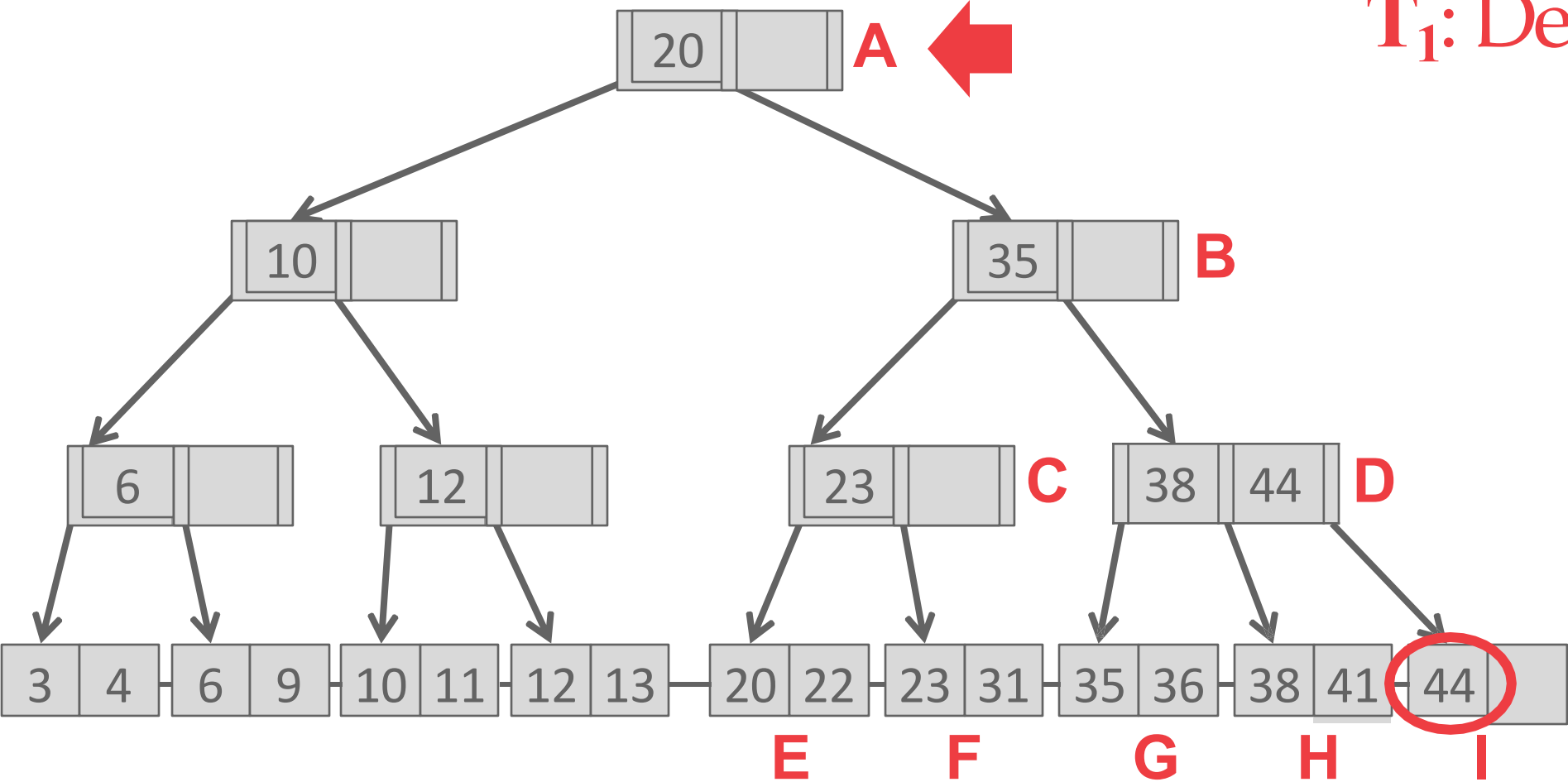
B+TREE MULTI-THREADED EXAMPLE

T₁: Delete 44



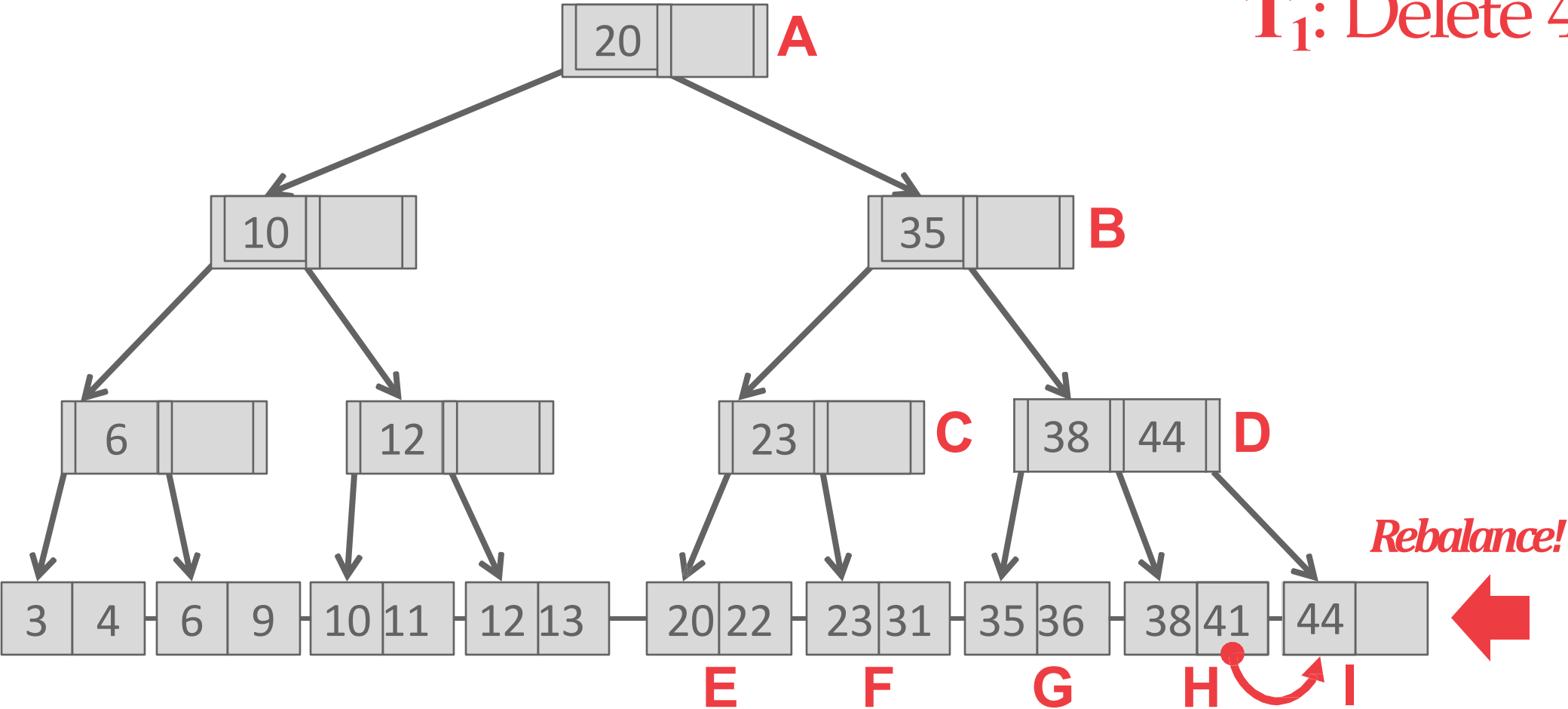
B+TREE MULTI-THREADED EXAMPLE

T₁: Delete 44

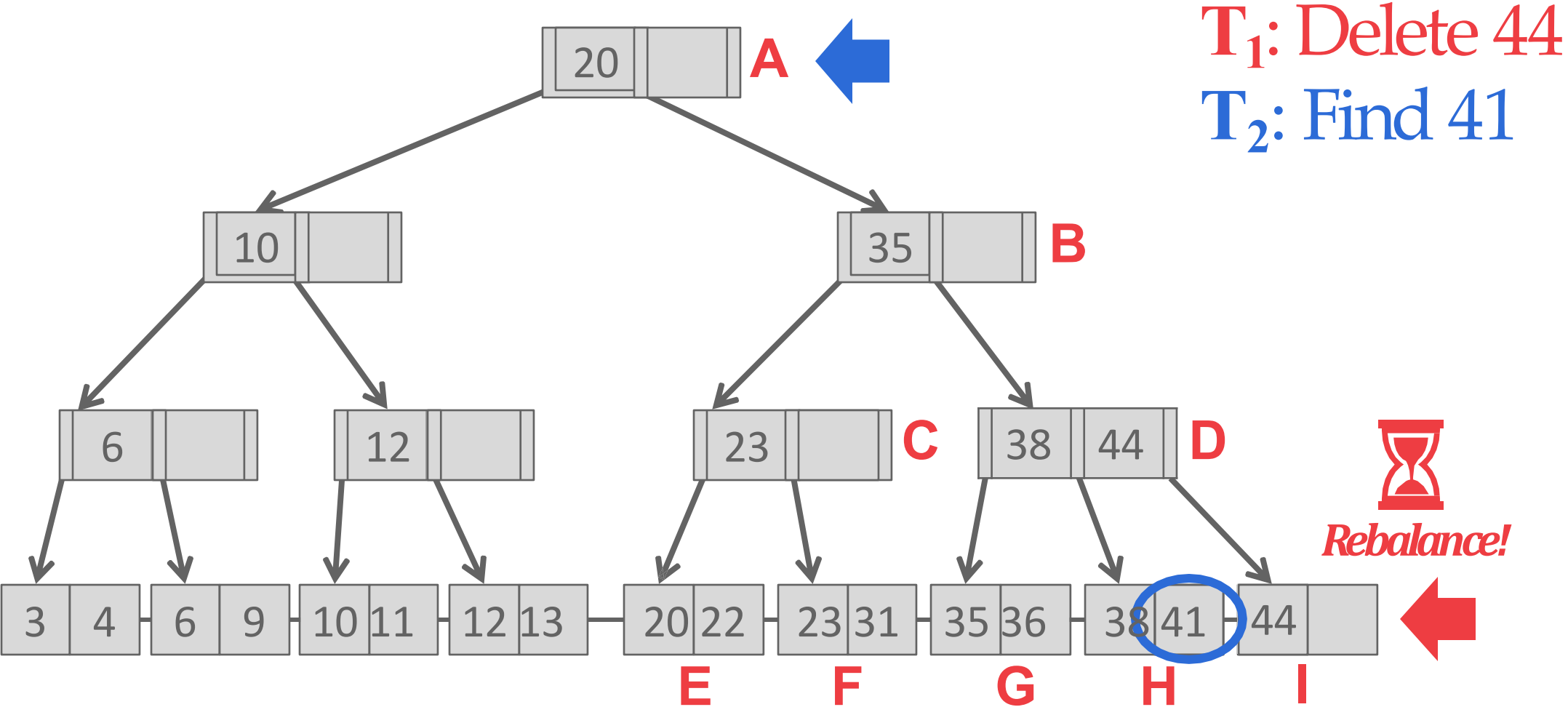


B+TREE MULTI-THREADED EXAMPLE

T₁: Delete 44



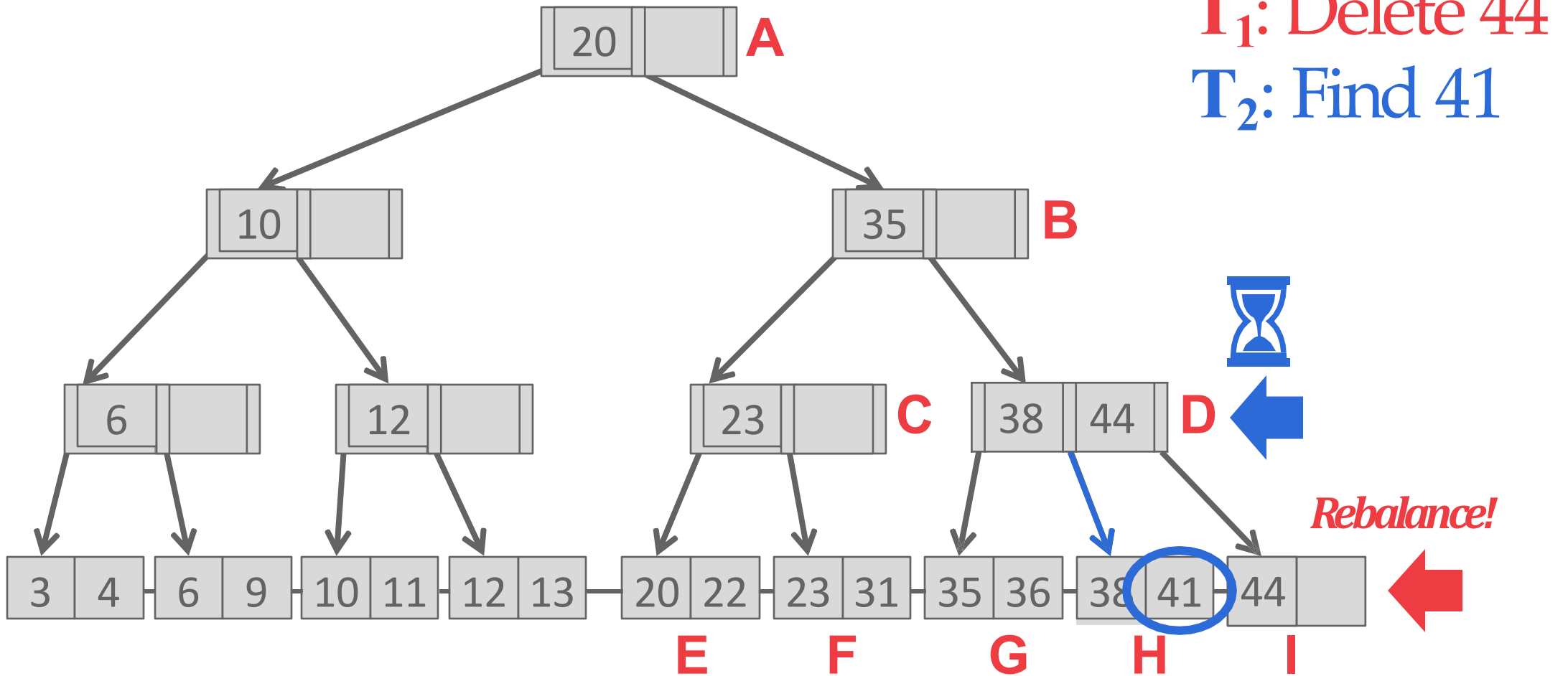
B+TREE MULTI-THREADED EXAMPLE



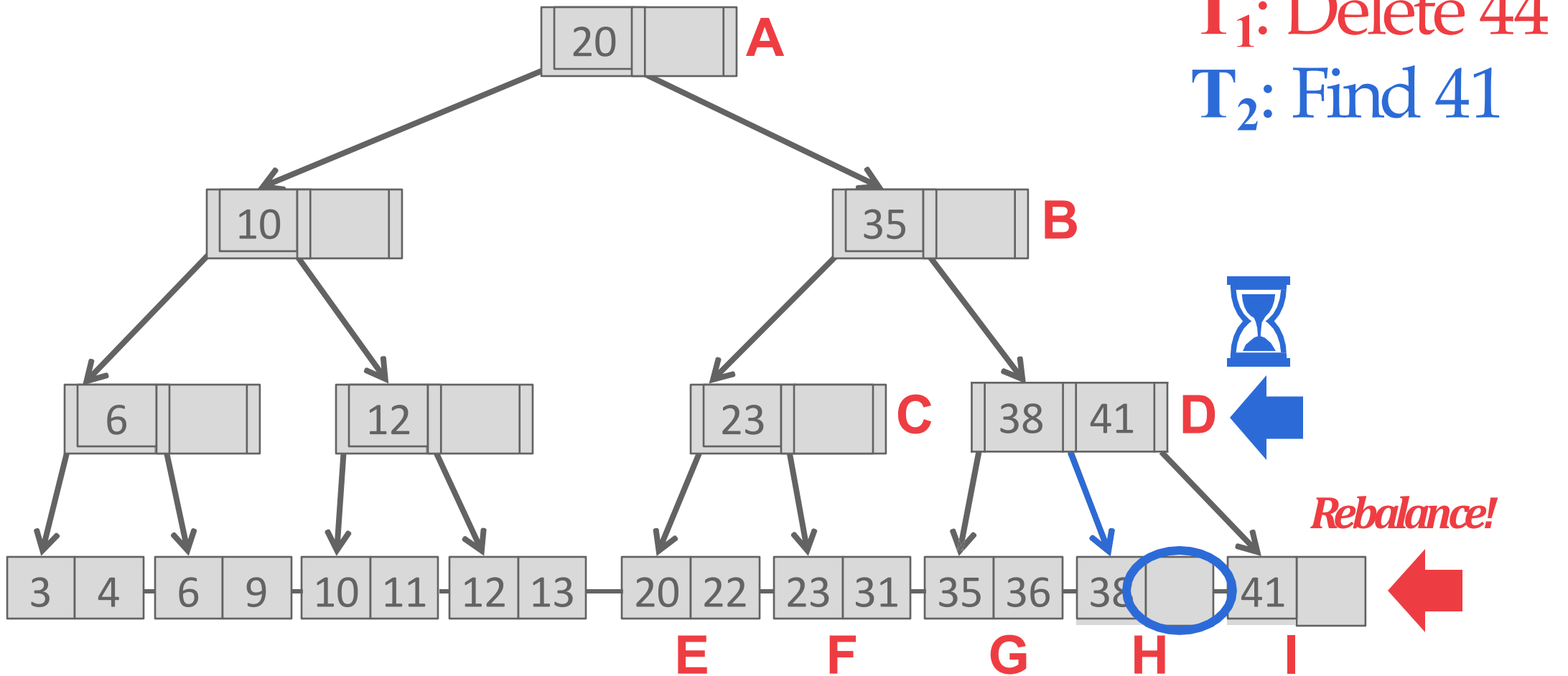
T₁: Delete 44
T₂: Find 41


Rebalance!

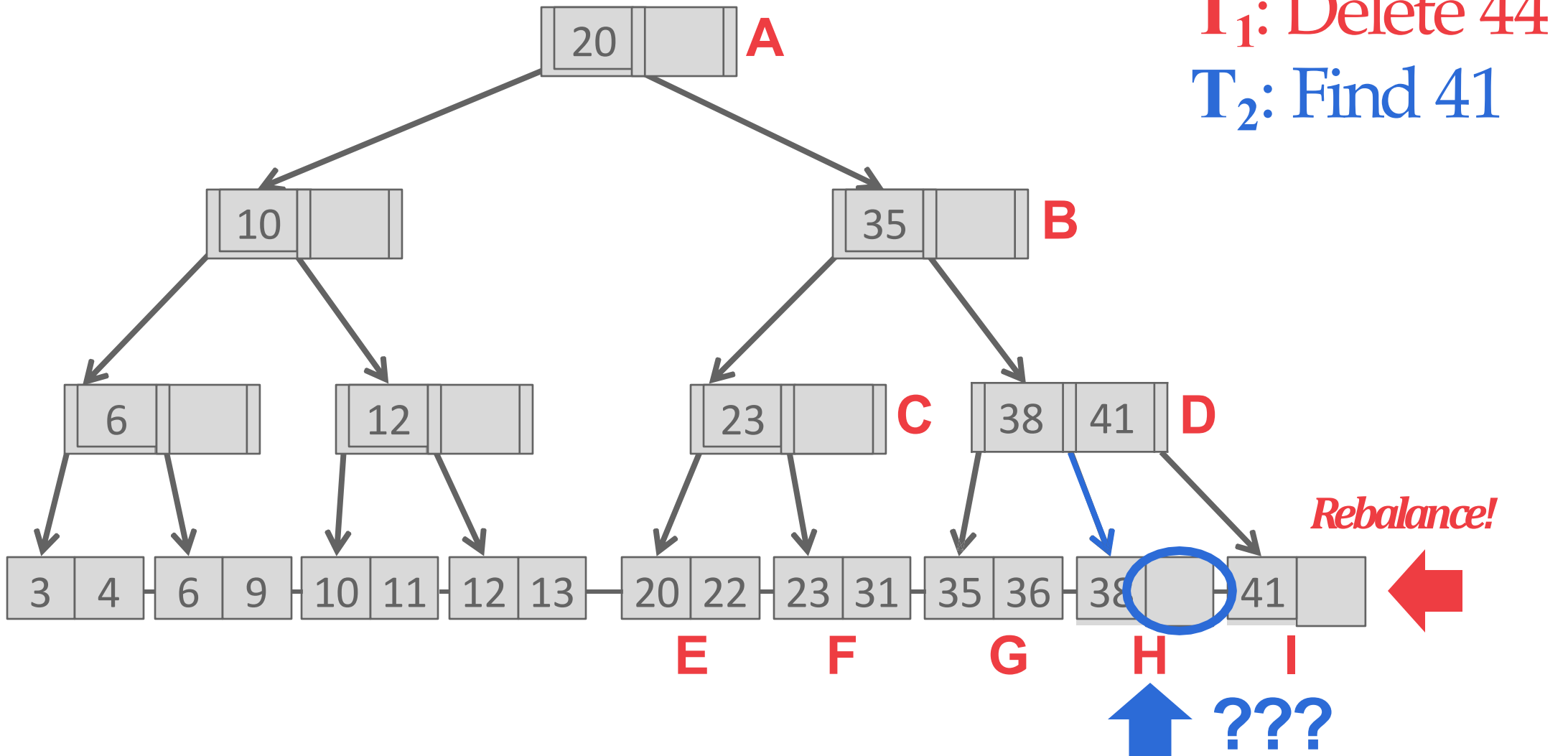
B+TREE MULTI-THREADED EXAMPLE



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.

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.

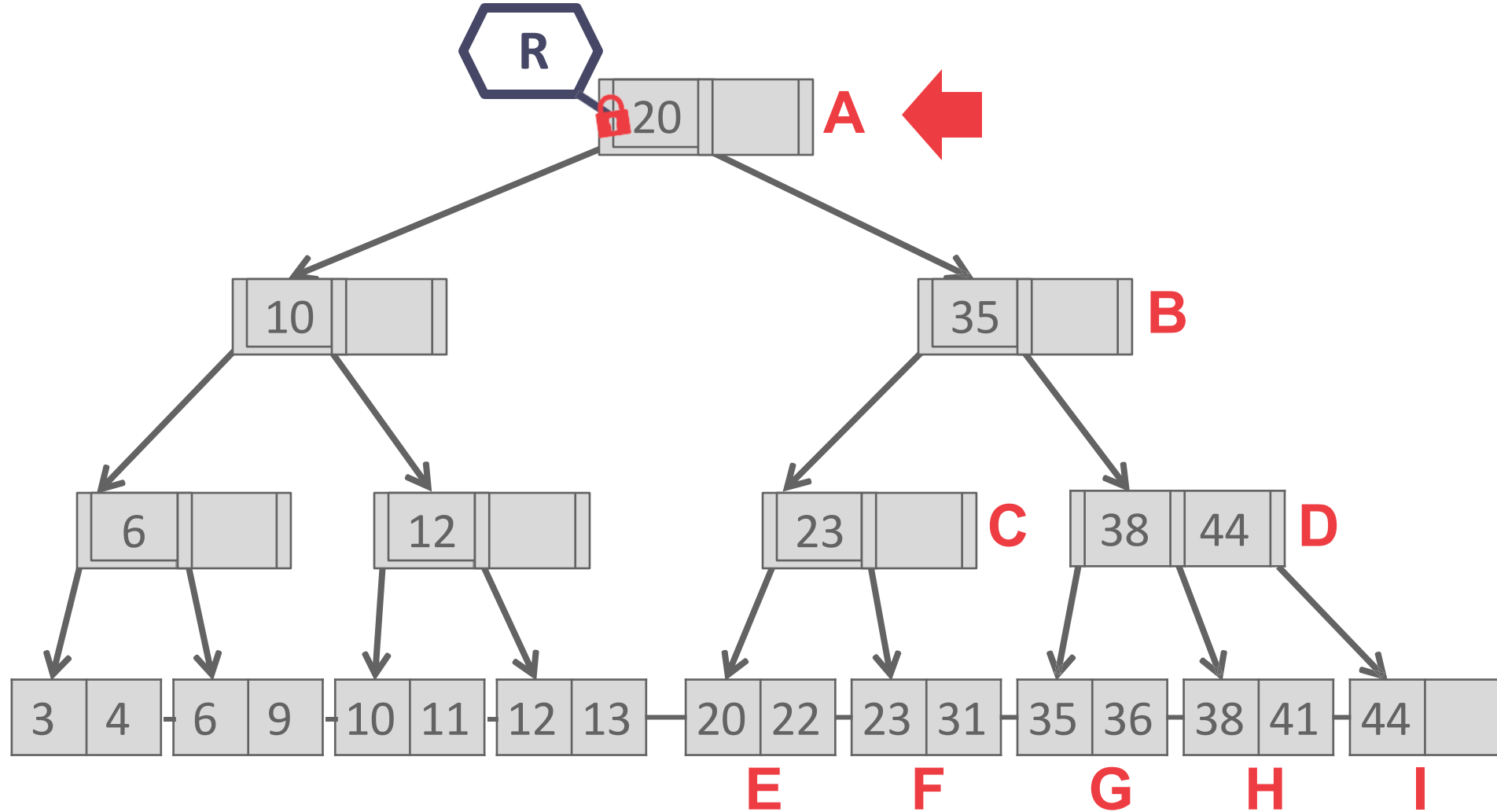
- Not full (on insertion)
- More than half-full (on deletion)

LATCH CRABBING/COUPLING

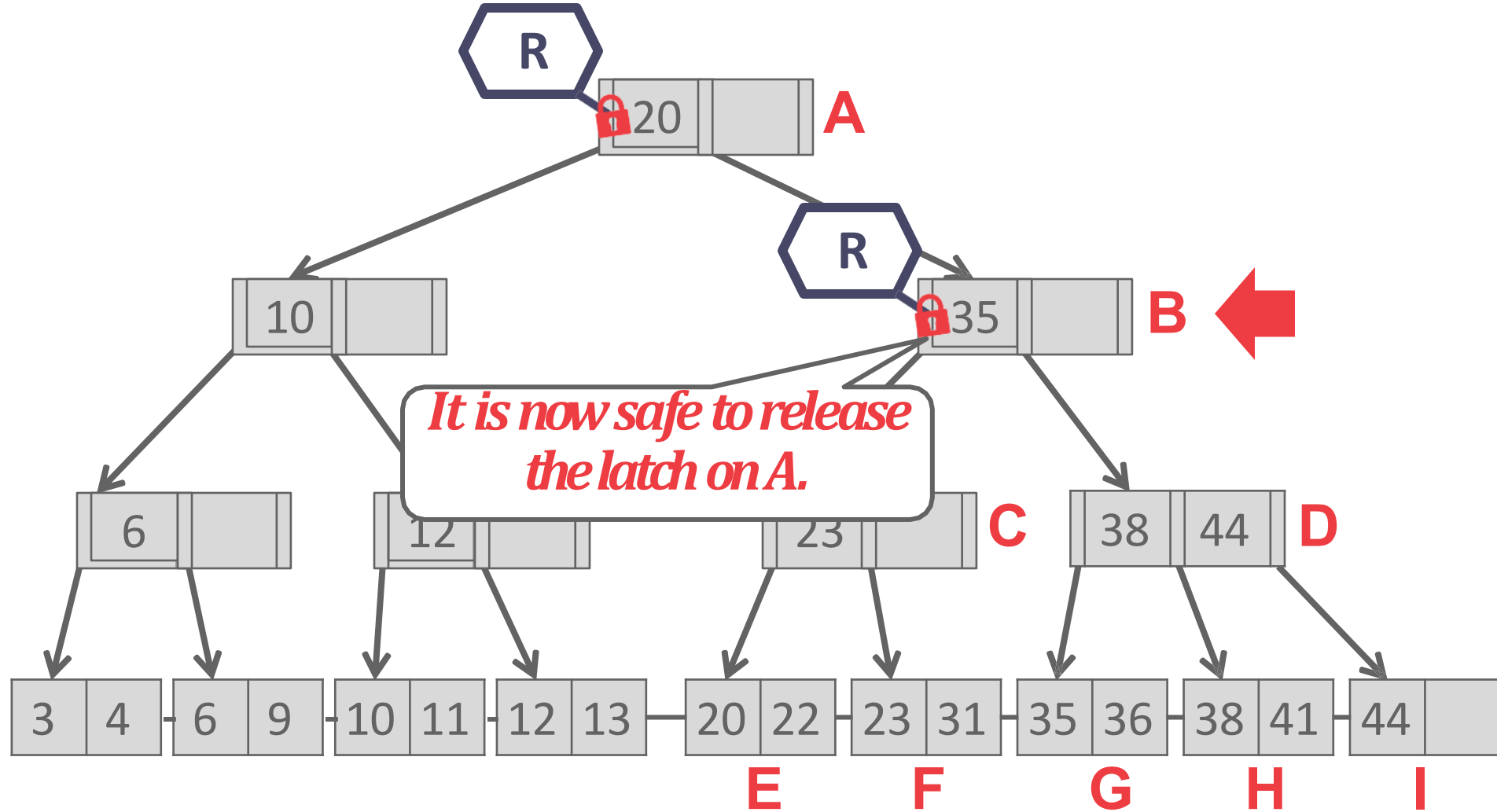
Find: Start at root and go down; repeatedly,
→ Acquire **R**latch on child
→ Then unlatch parent

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.

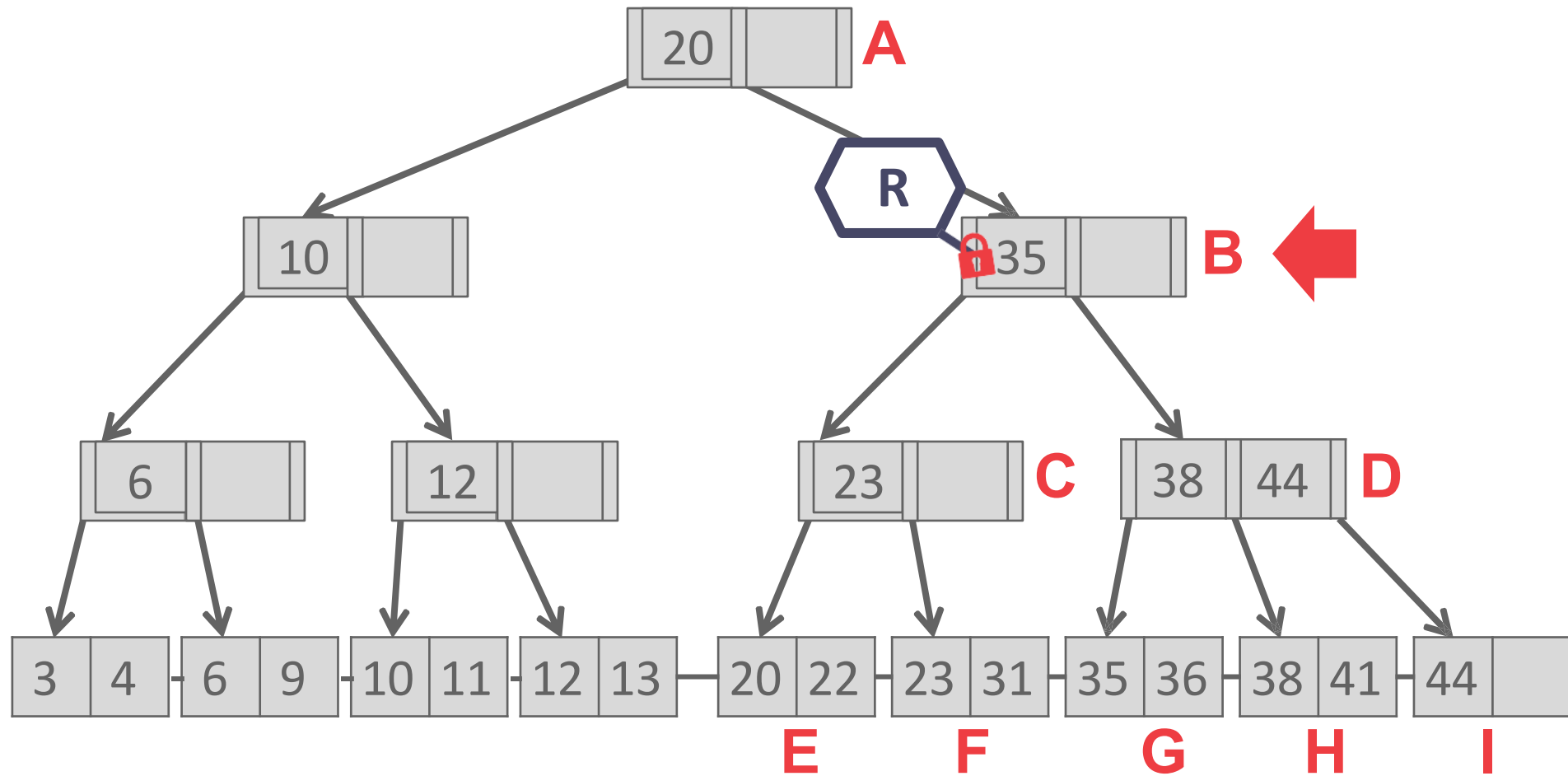
EXAMPLE #1 – FIND 38



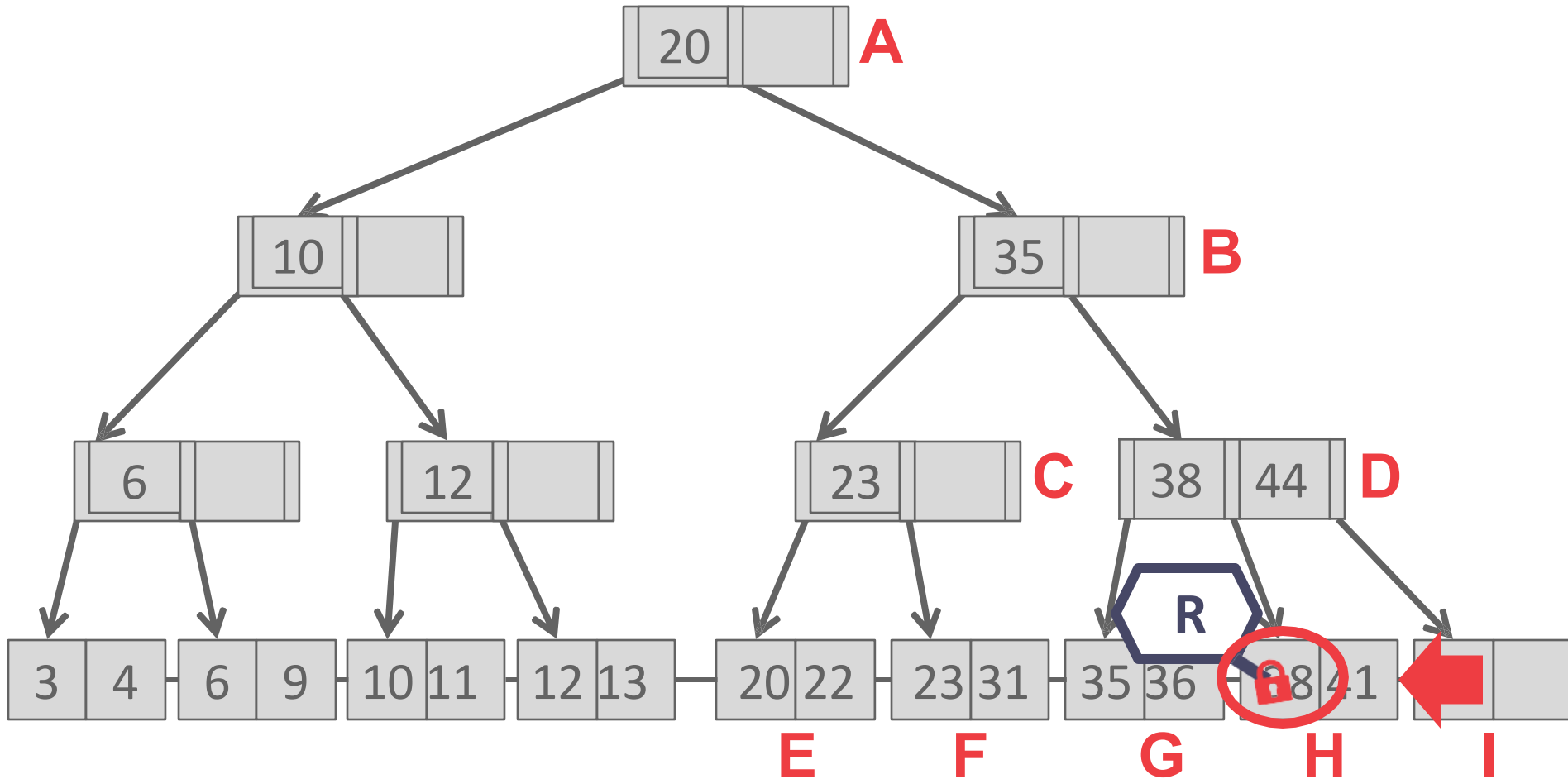
EXAMPLE #1 – FIND 38



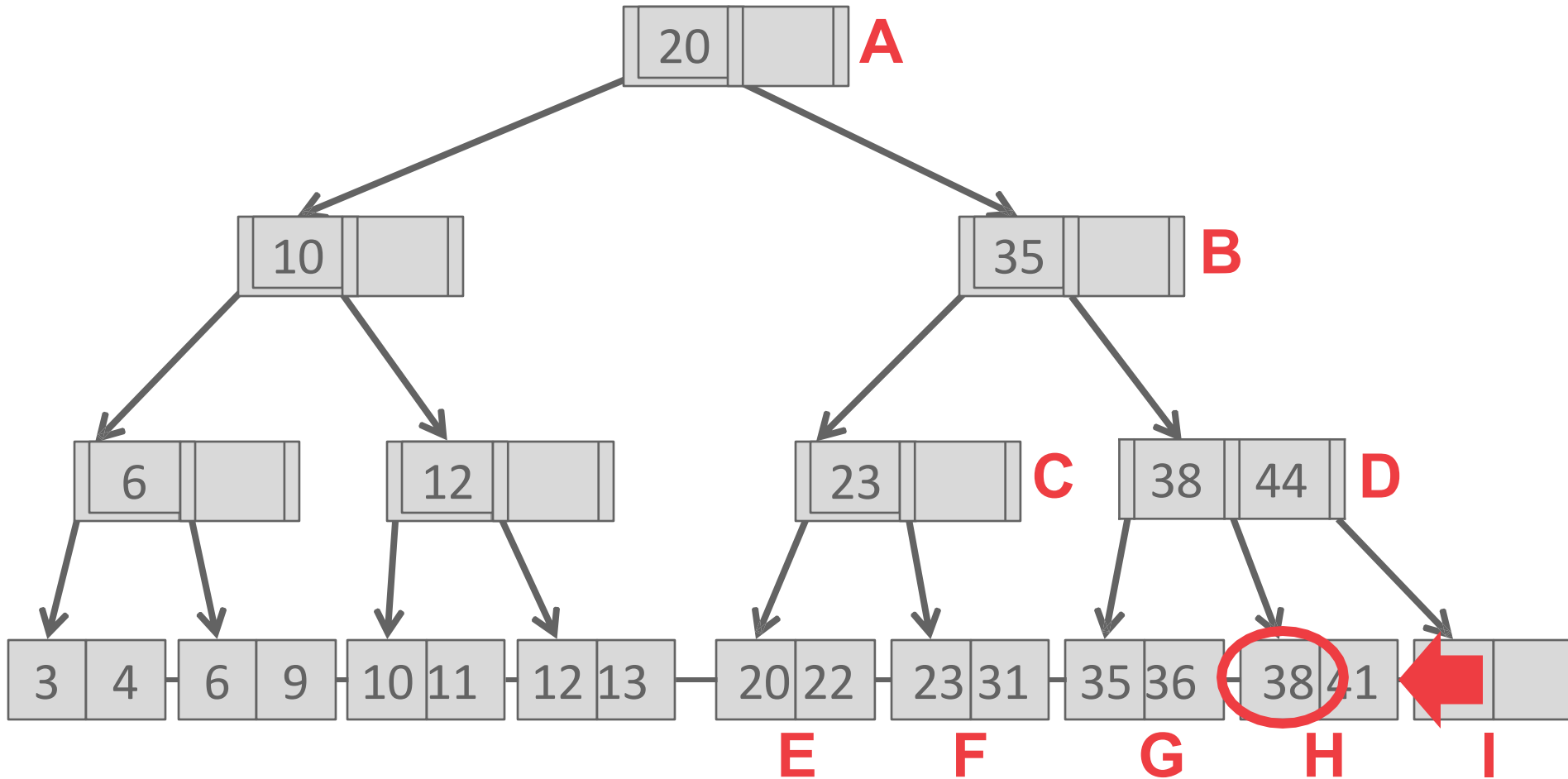
EXAMPLE #1 – FIND 38



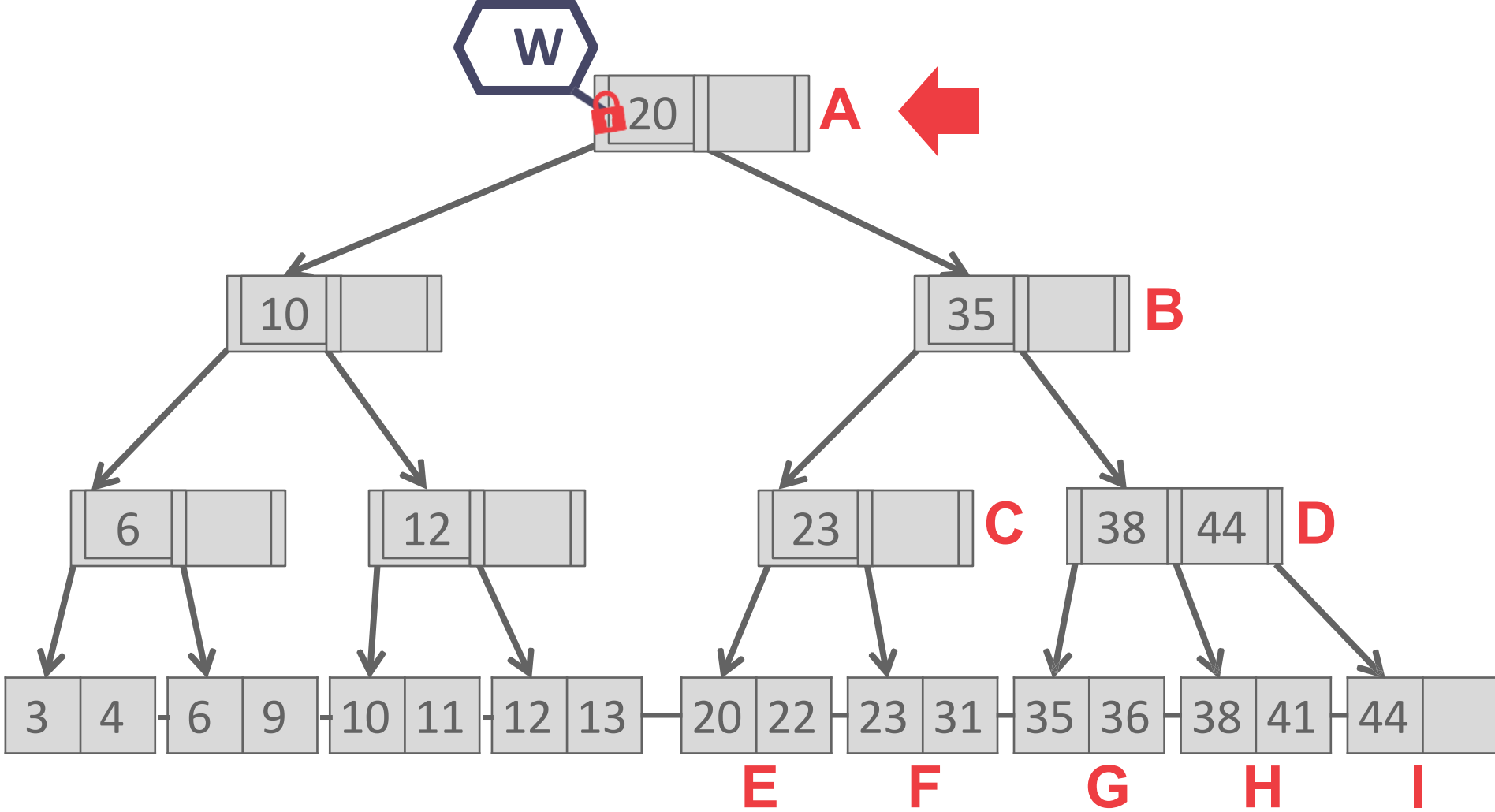
EXAMPLE #1 – FIND 38



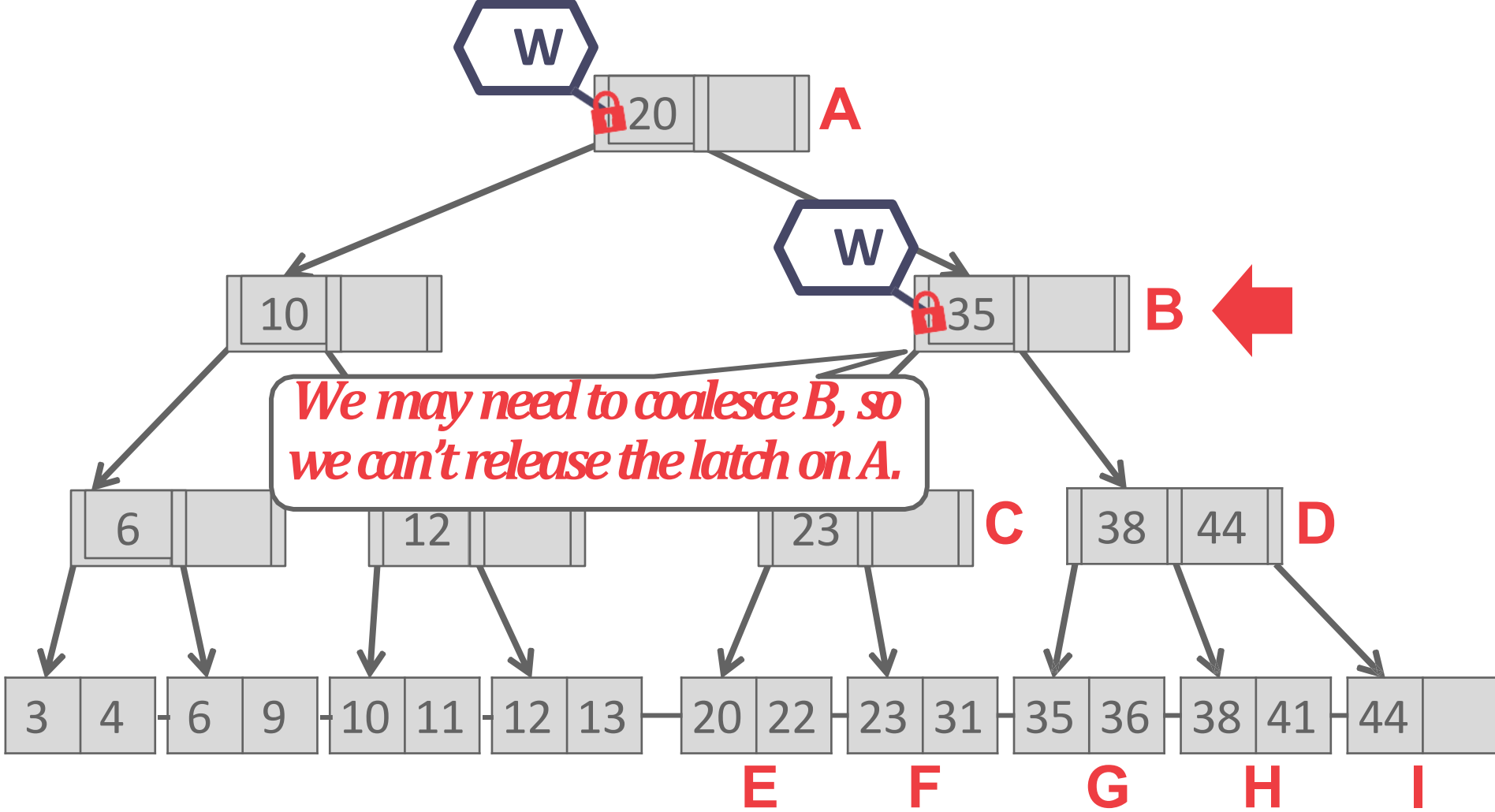
EXAMPLE #1 – FIND 38



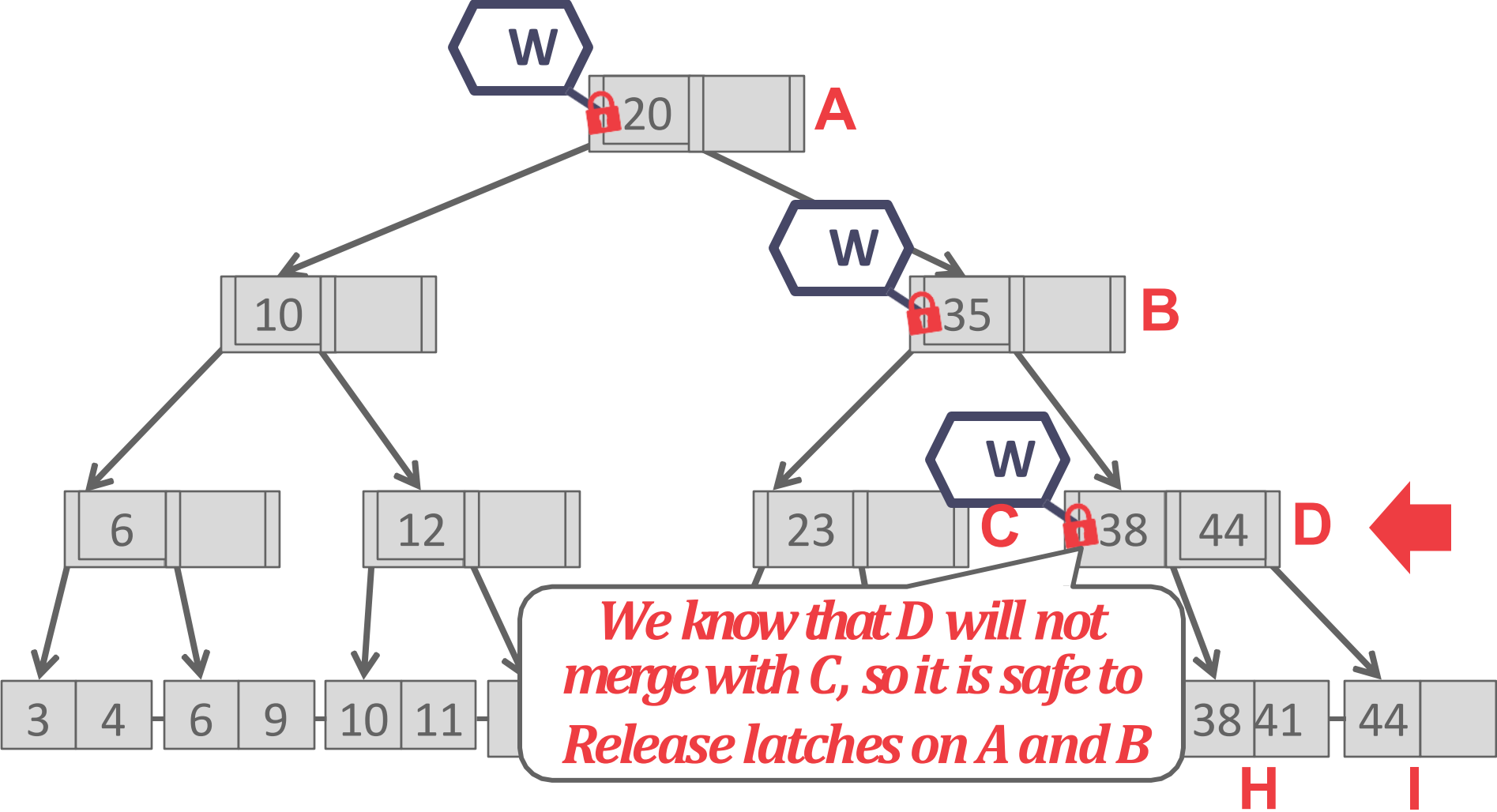
EXAMPLE #2 – DELETE 38



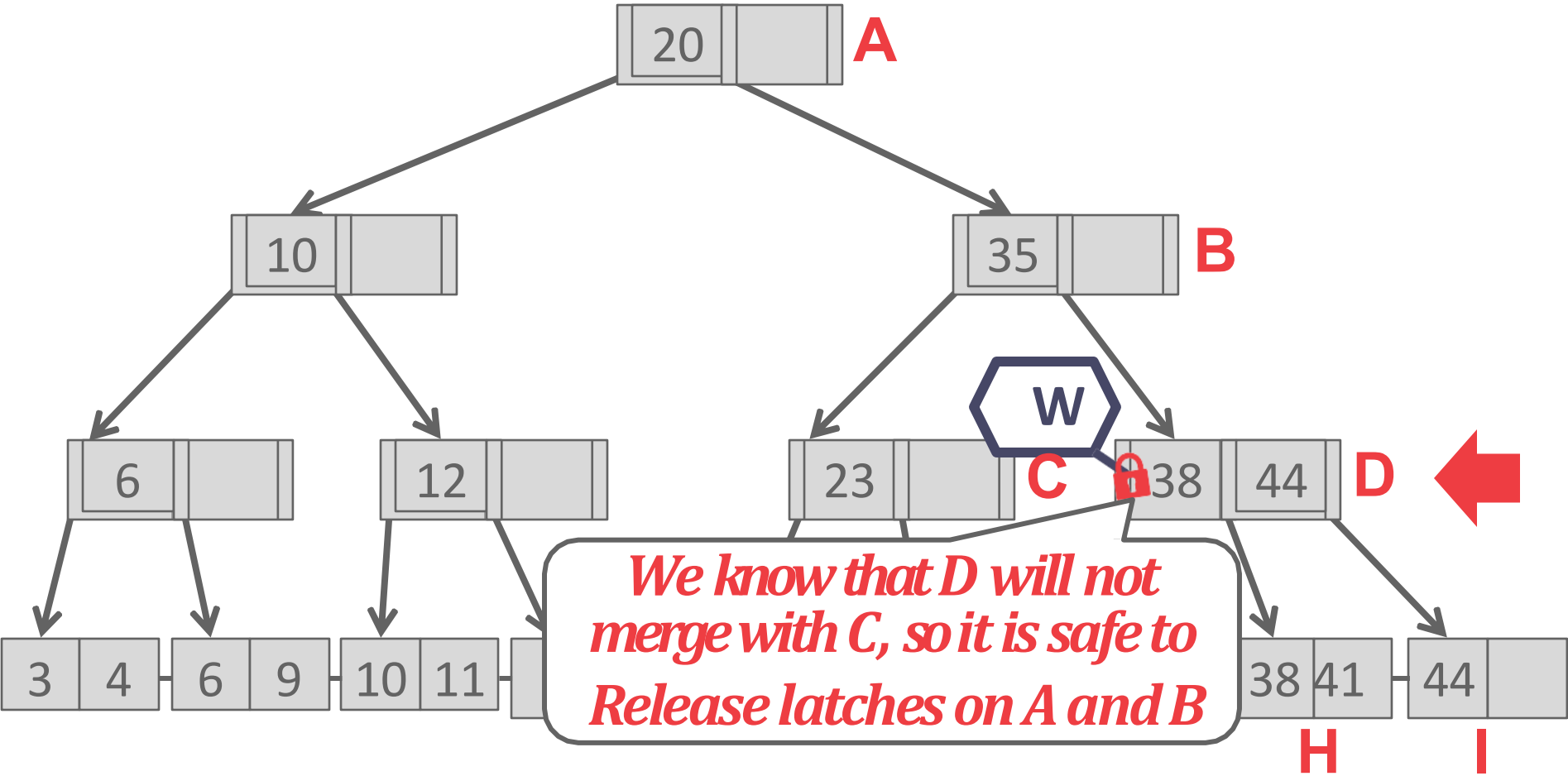
EXAMPLE #2 – DELETE 38



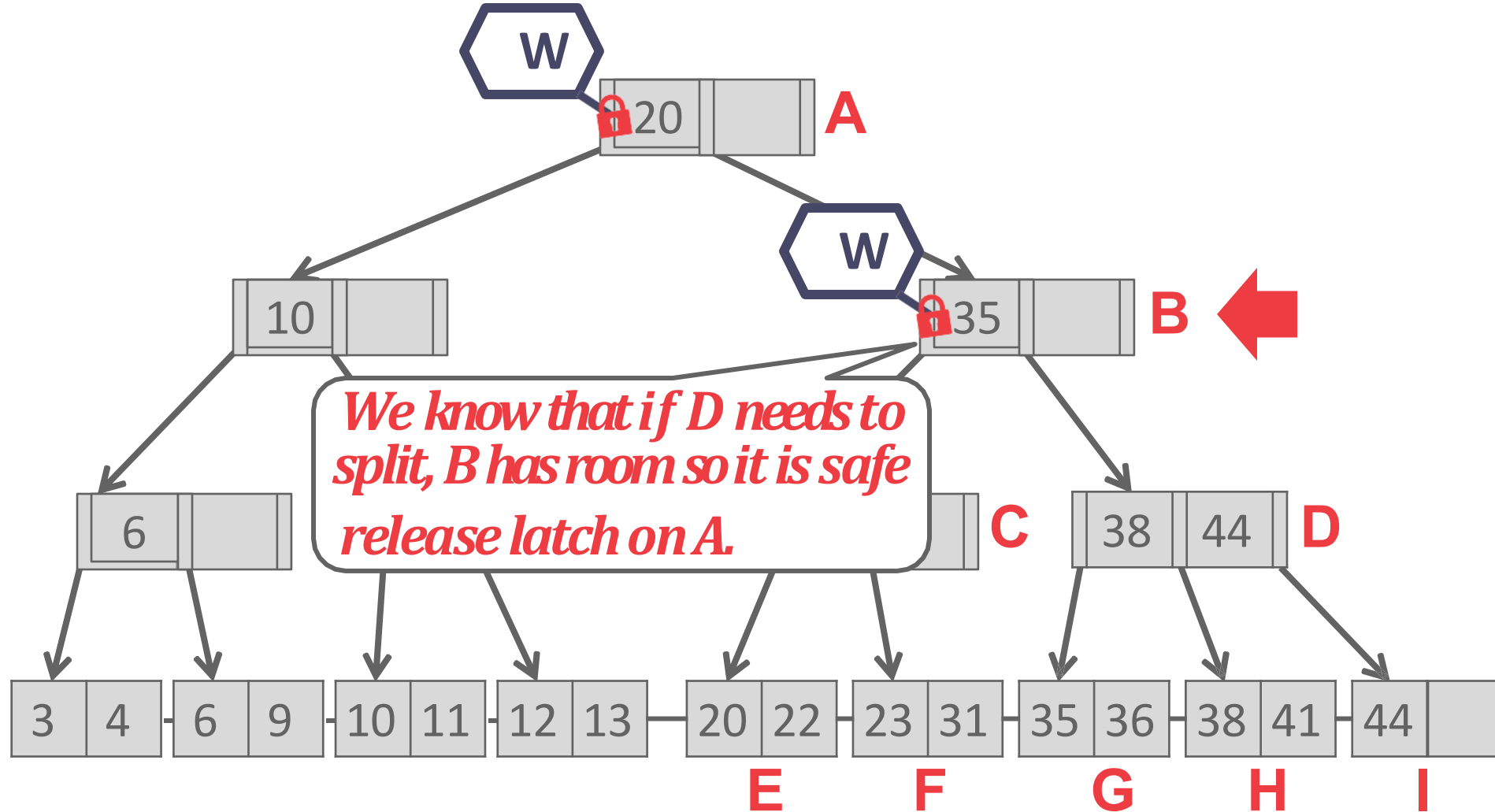
EXAMPLE #2 – DELETE 38



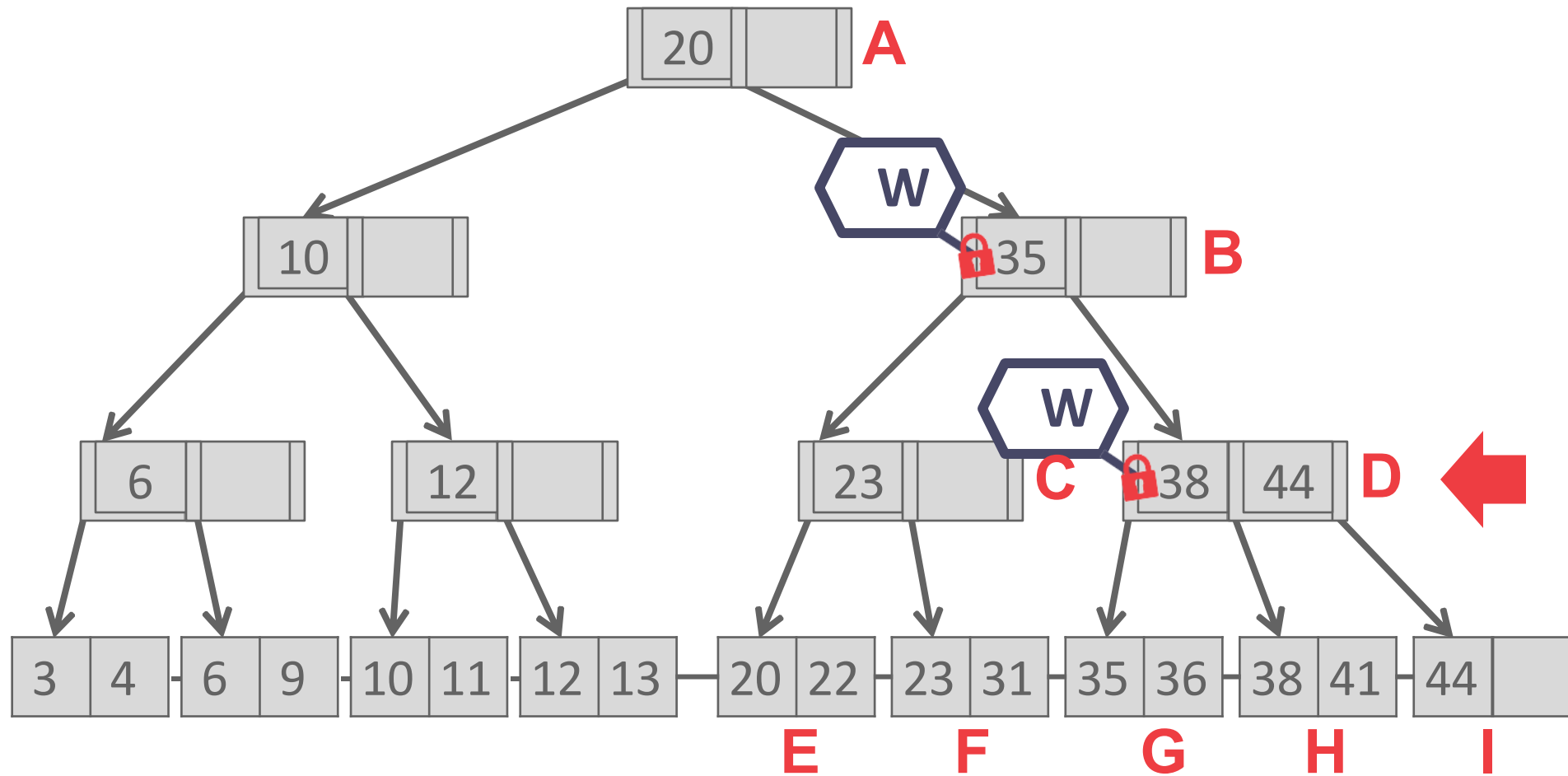
EXAMPLE #2 – DELETE 38



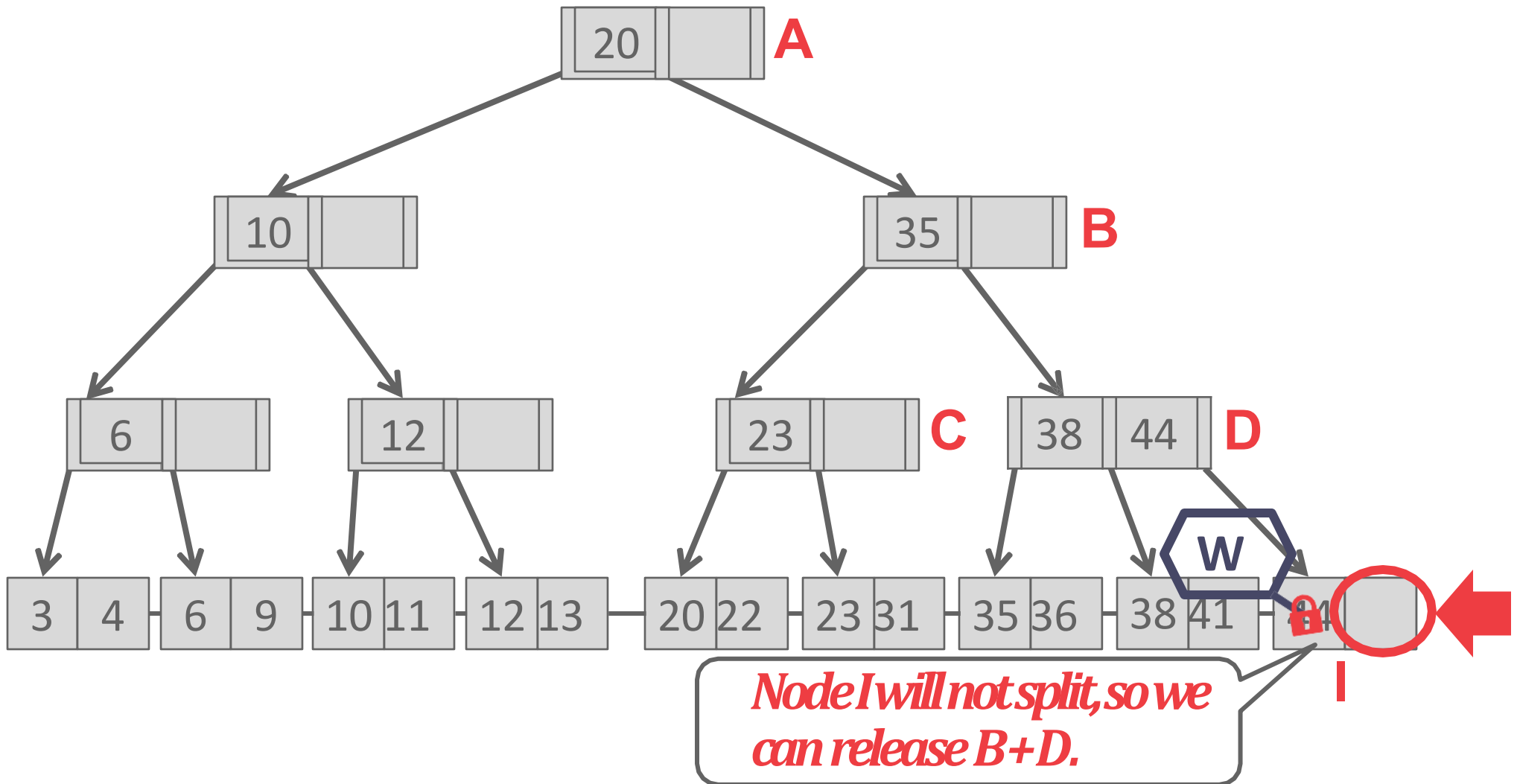
EXAMPLE #3 – INSERT 45



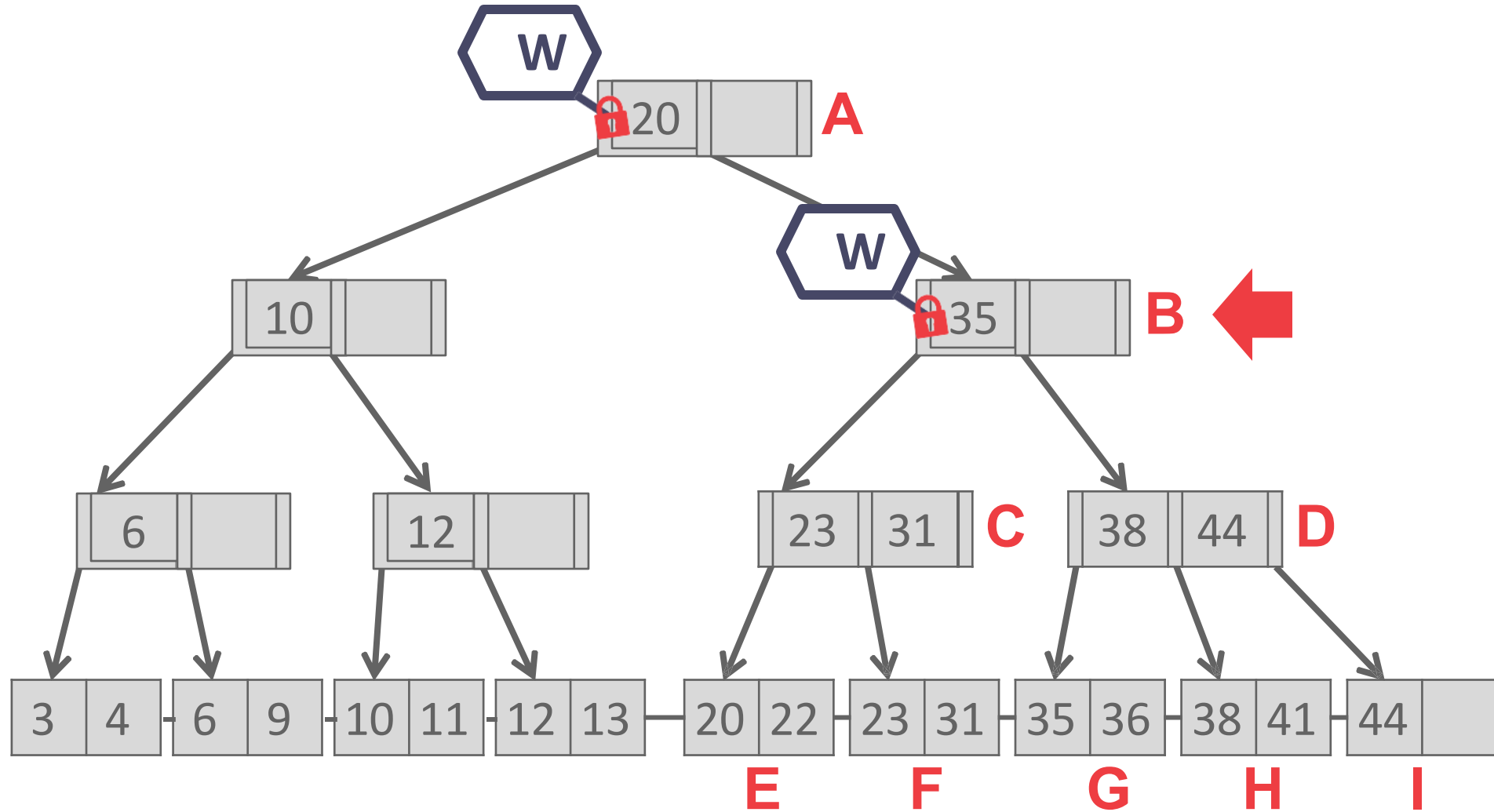
EXAMPLE #3 – INSERT 45



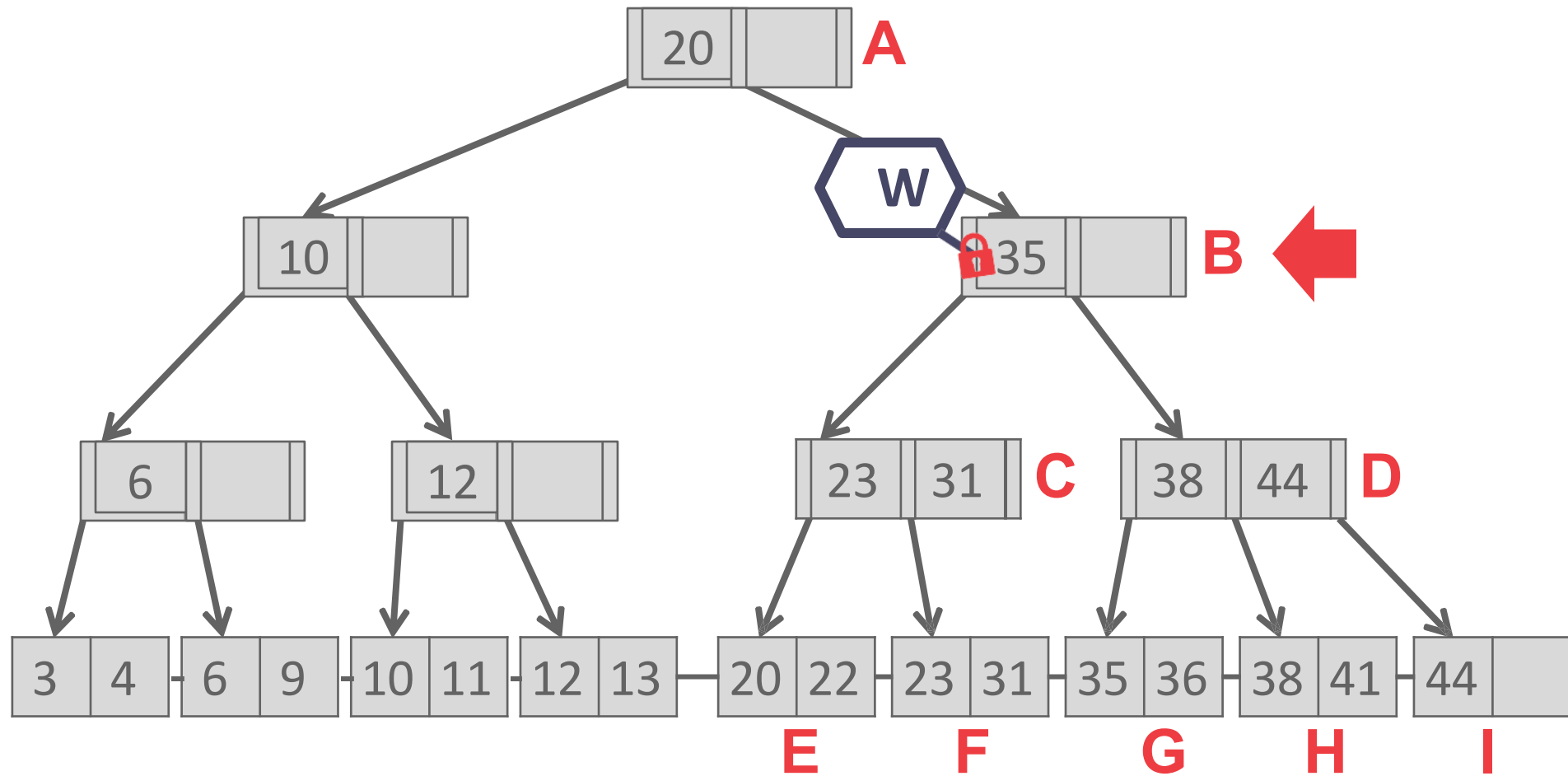
EXAMPLE #3 – INSERT 45



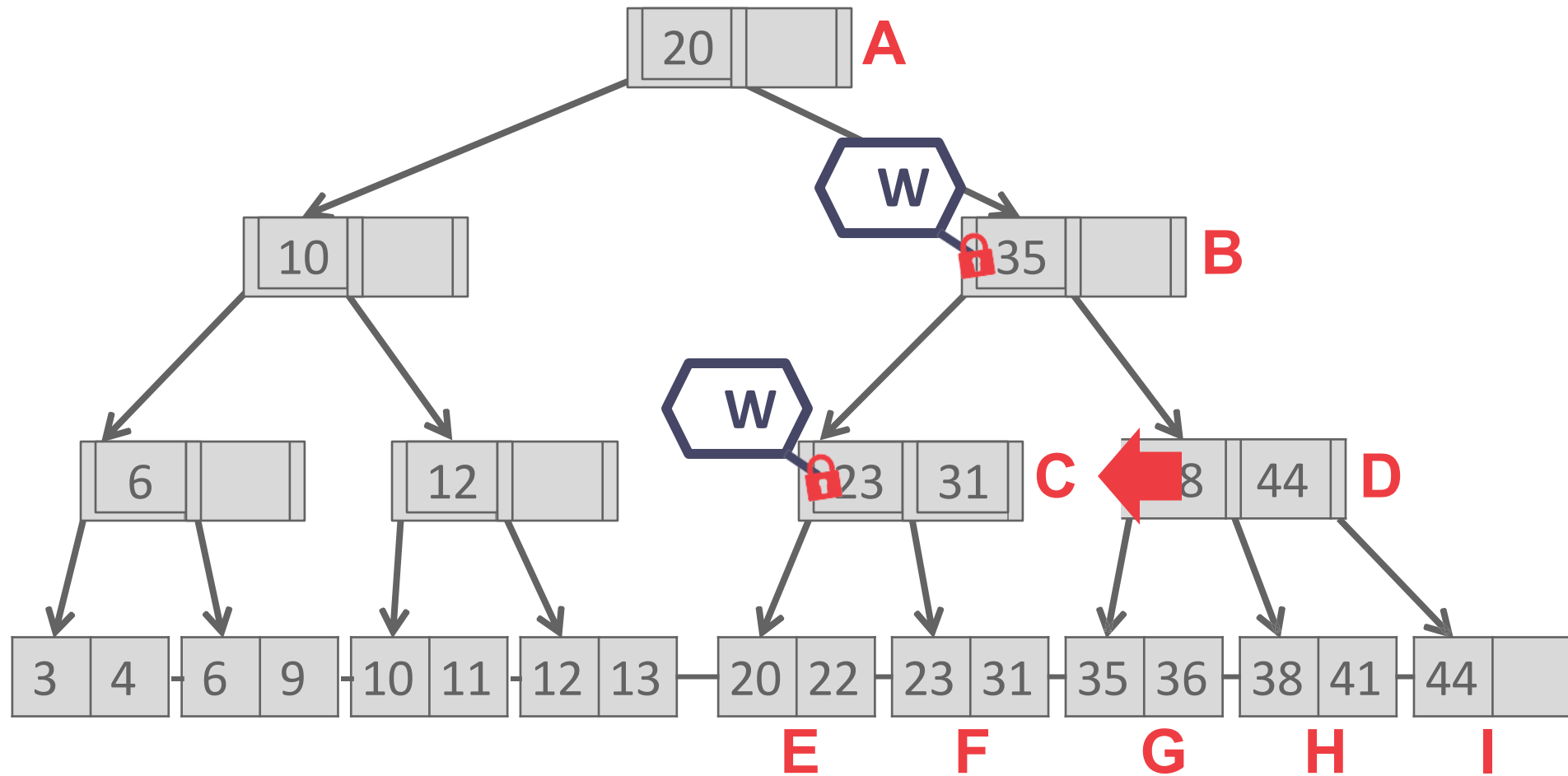
EXAMPLE #4 – INSERT 25



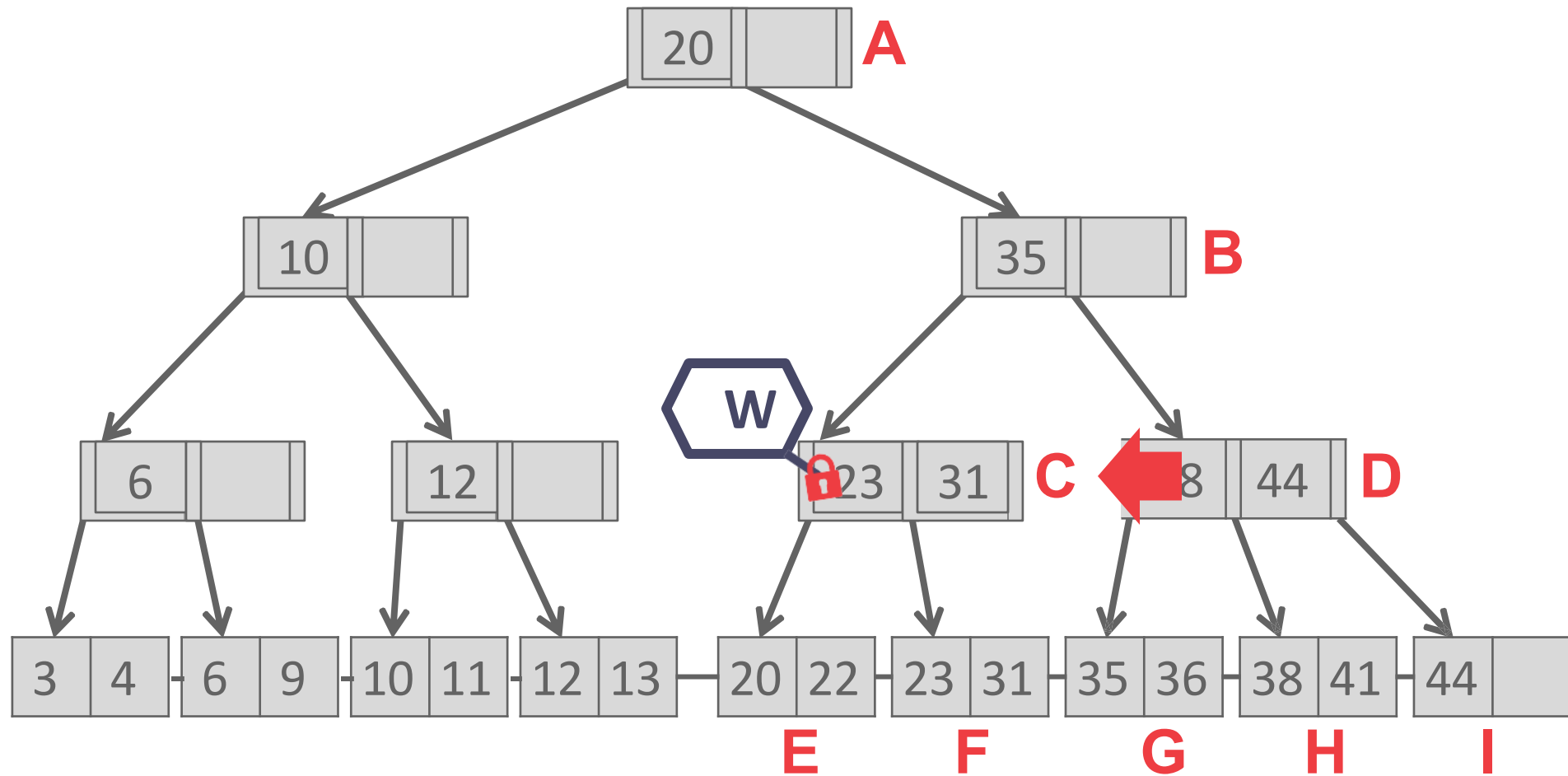
EXAMPLE #4 – INSERT 25



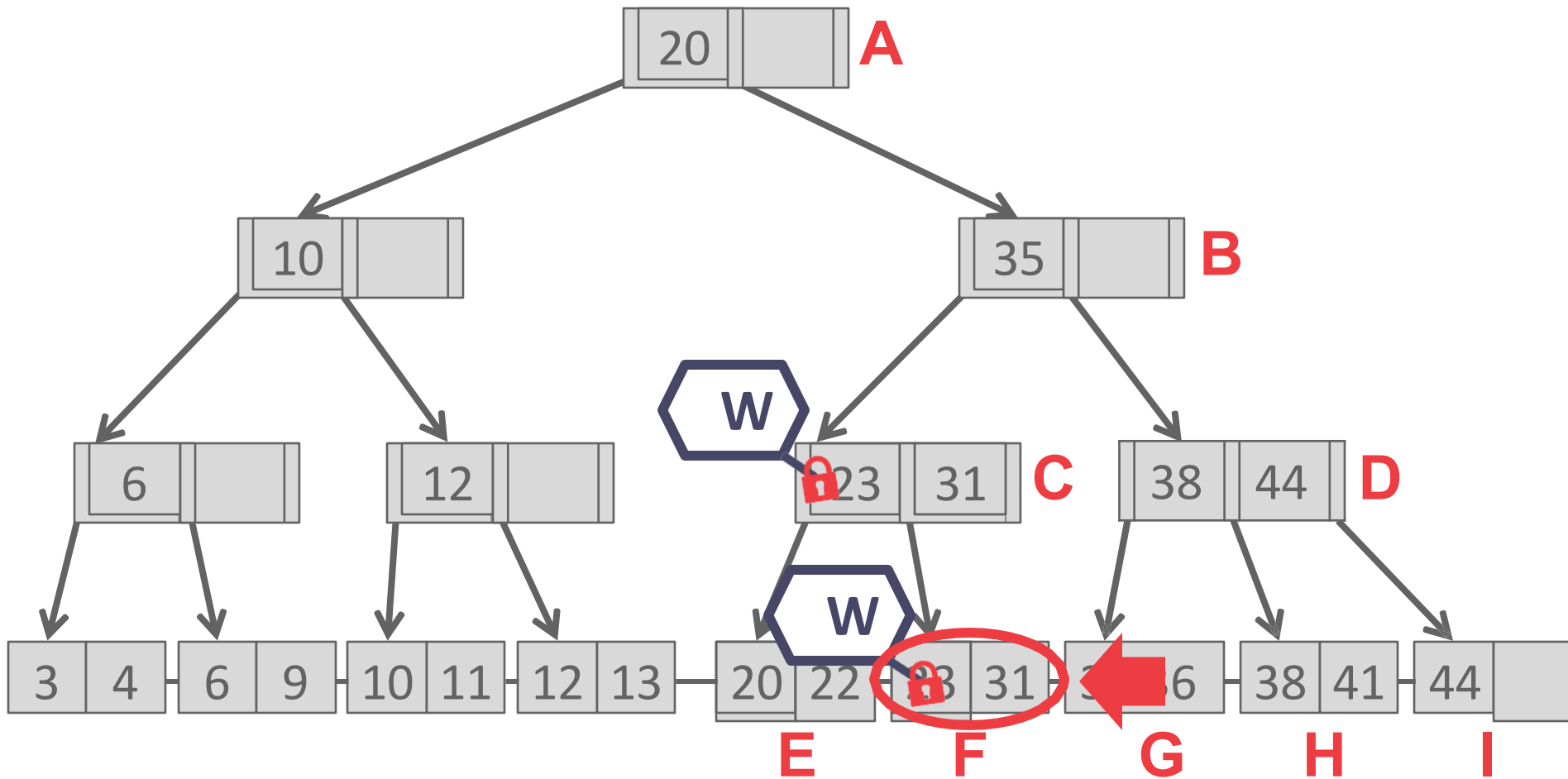
EXAMPLE #4 – INSERT 25



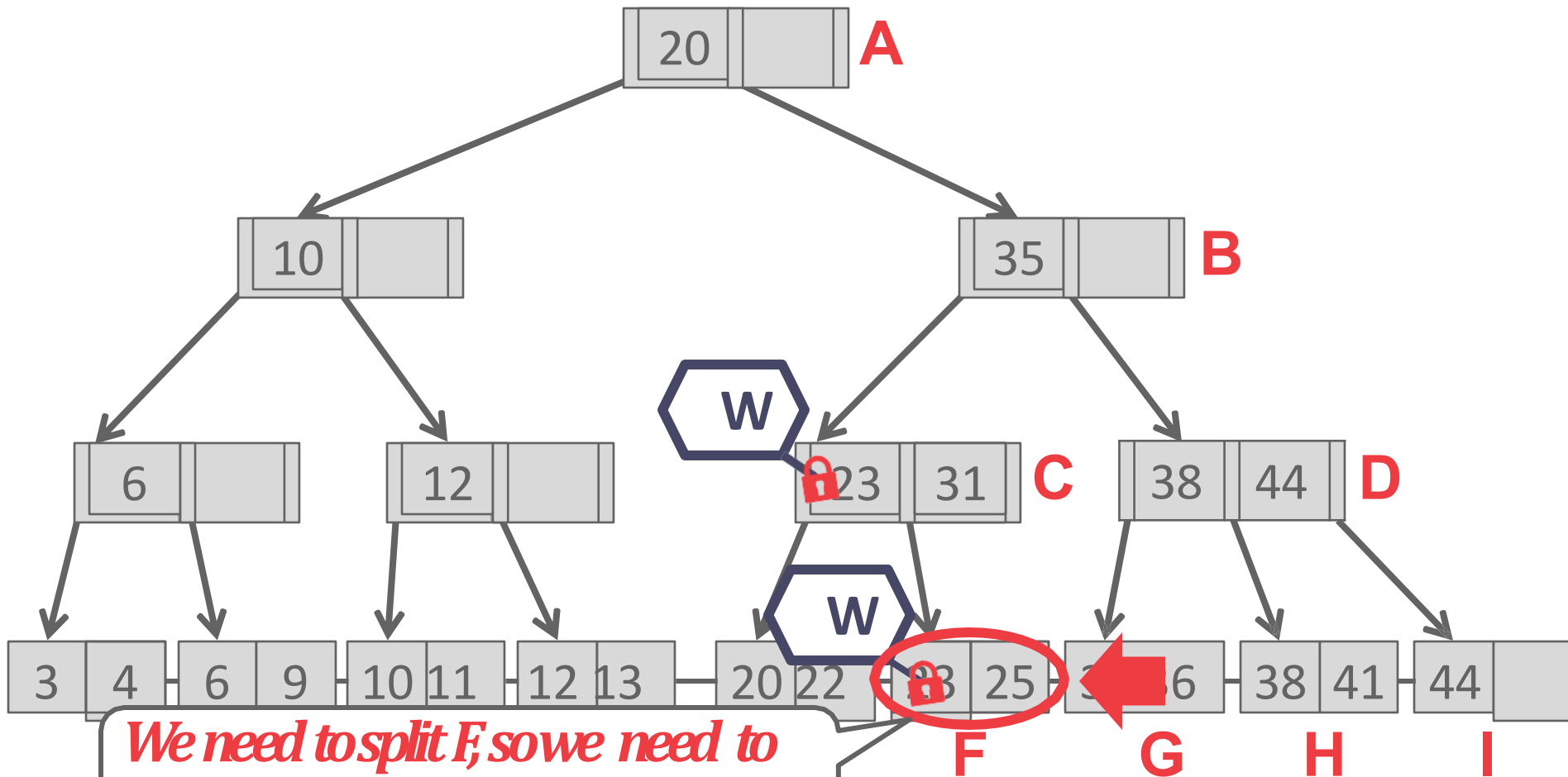
EXAMPLE #4 – INSERT 25



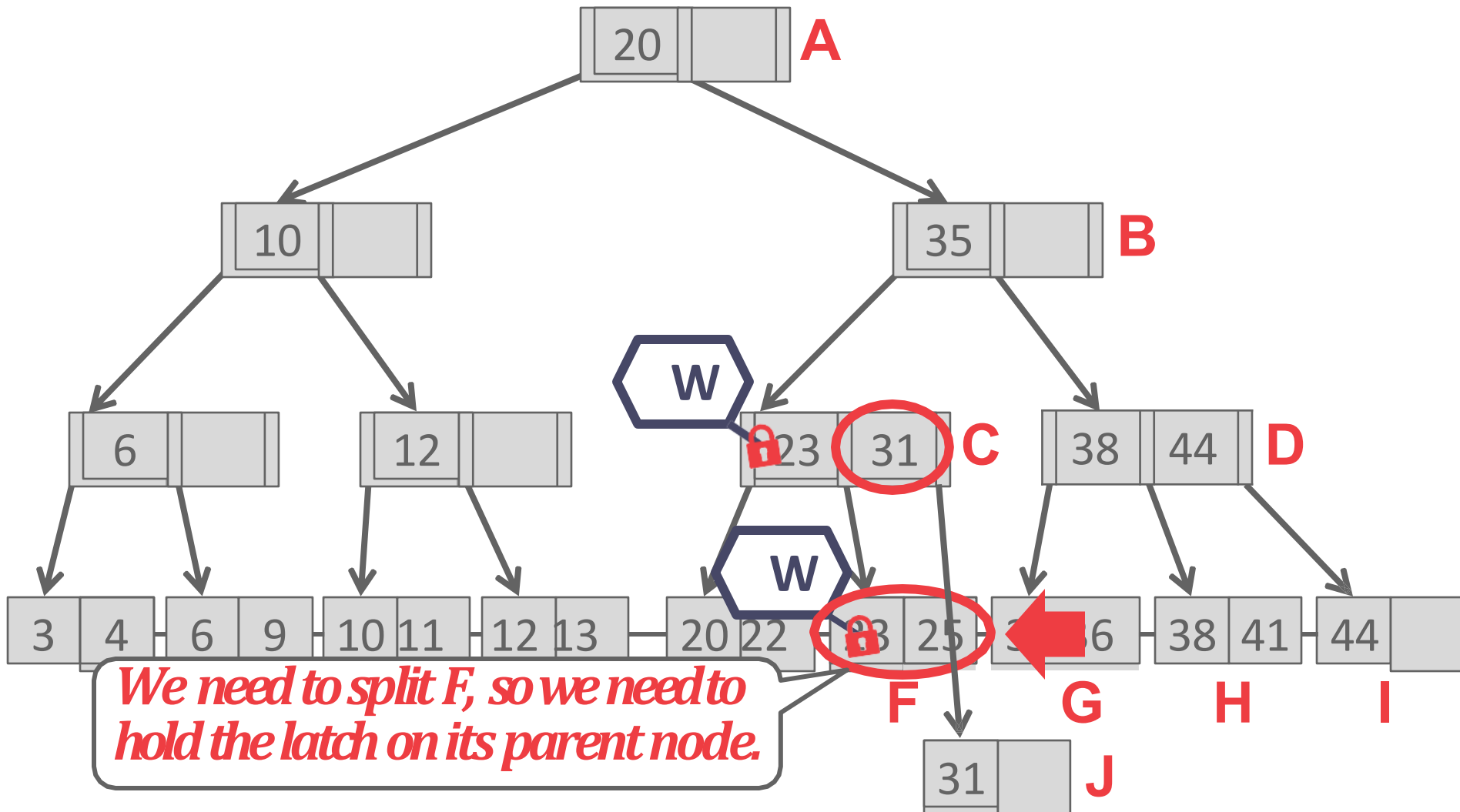
EXAMPLE #4 – INSERT 25



EXAMPLE #4 – INSERT 25

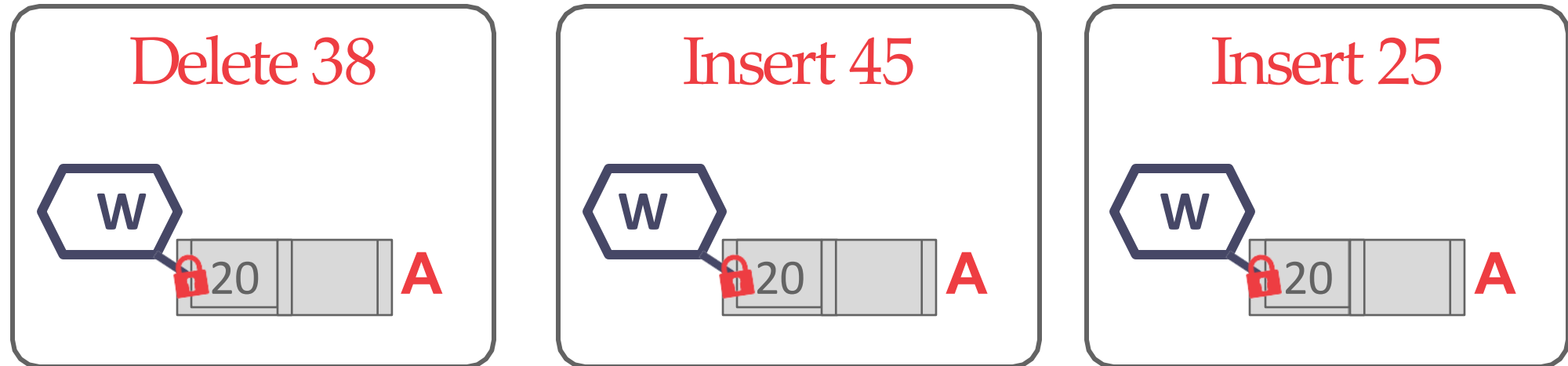


EXAMPLE #4 – INSERT 25



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

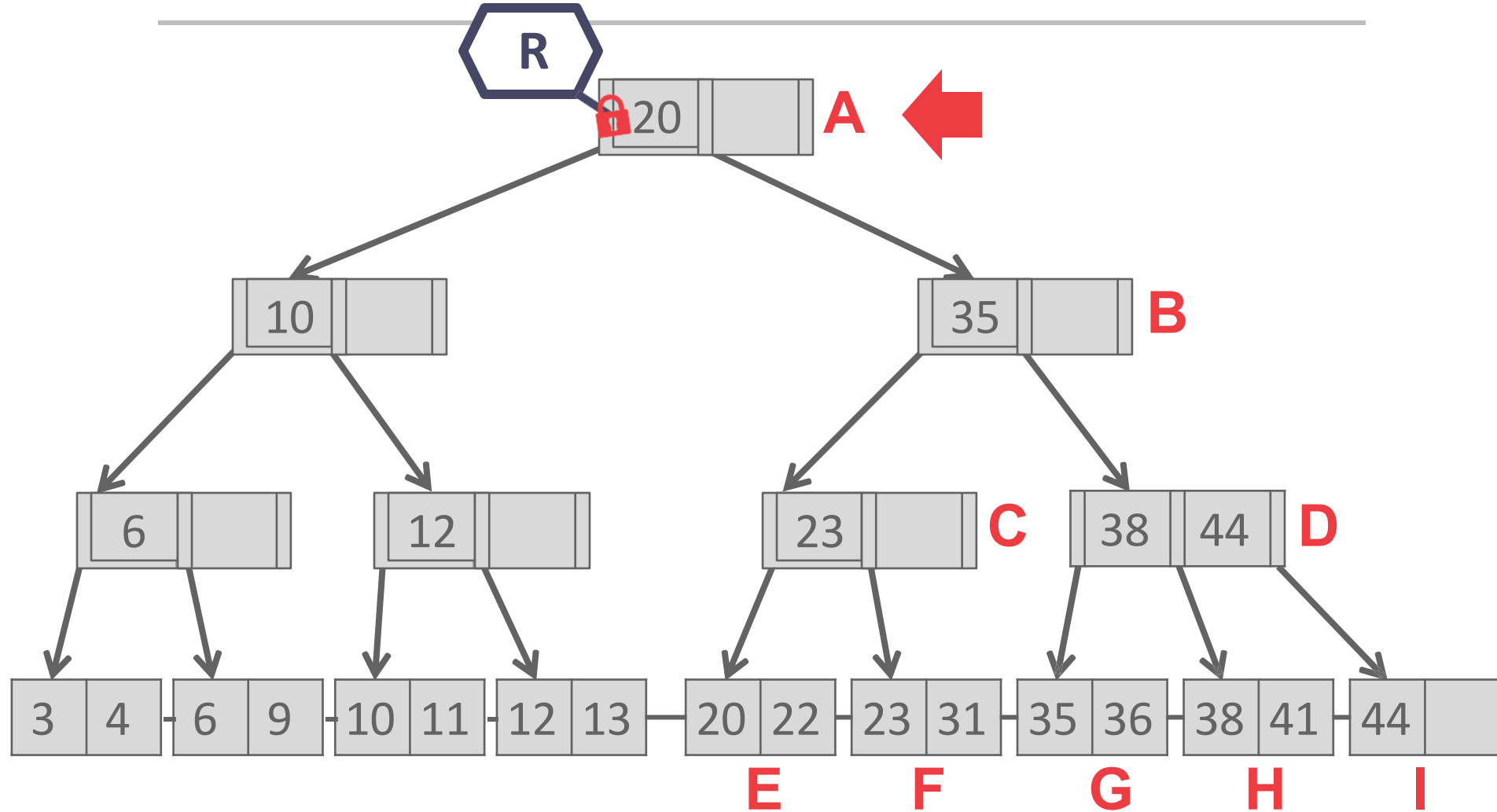
Search: Same as before.

Insert/Delete:

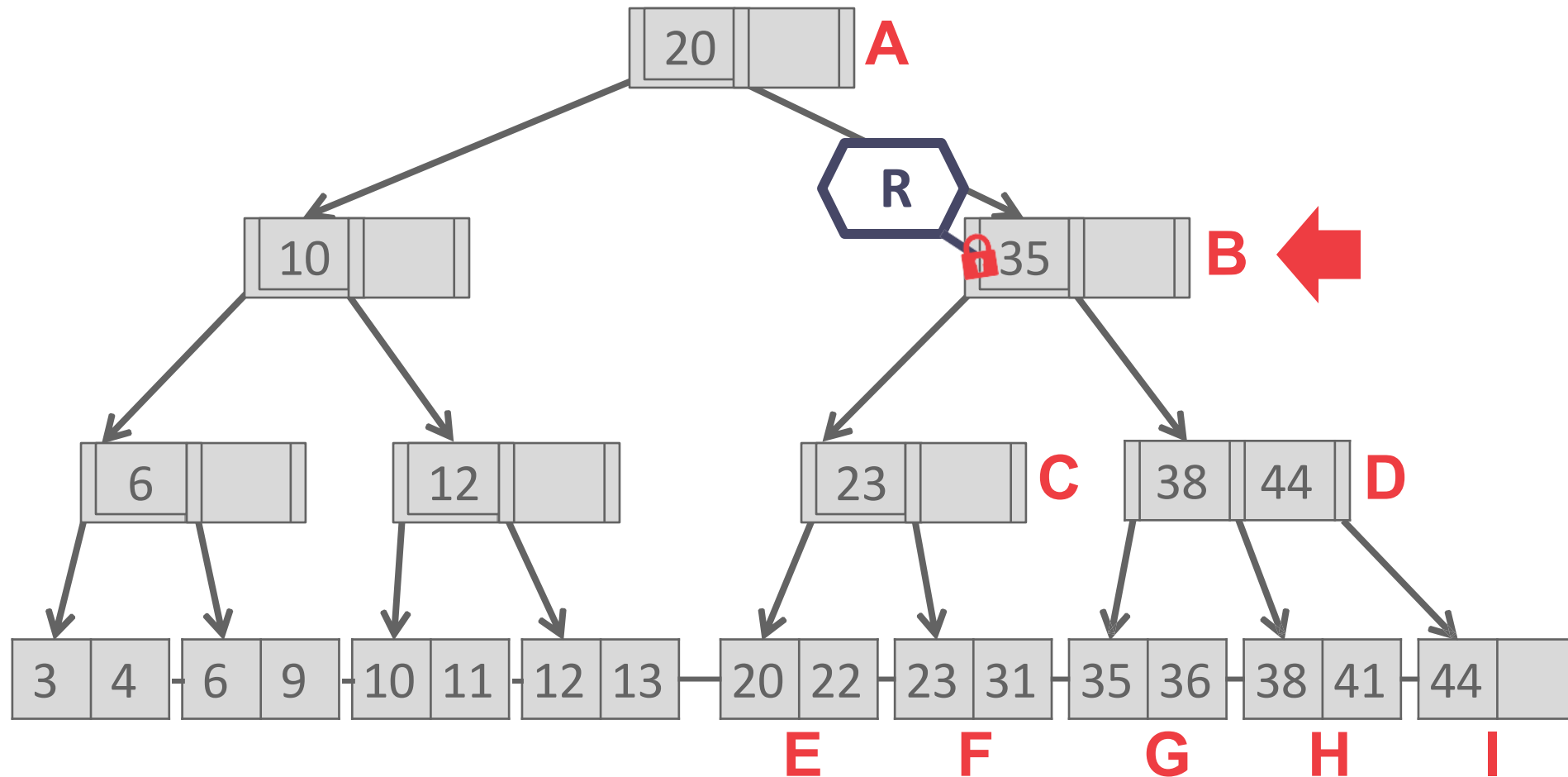
- Set latches as if for search, get to leaf, and set **W**latch on leaf.
- If leaf is not safe, release all latches, and restart thread using previous insert/delete protocol with write latches.

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

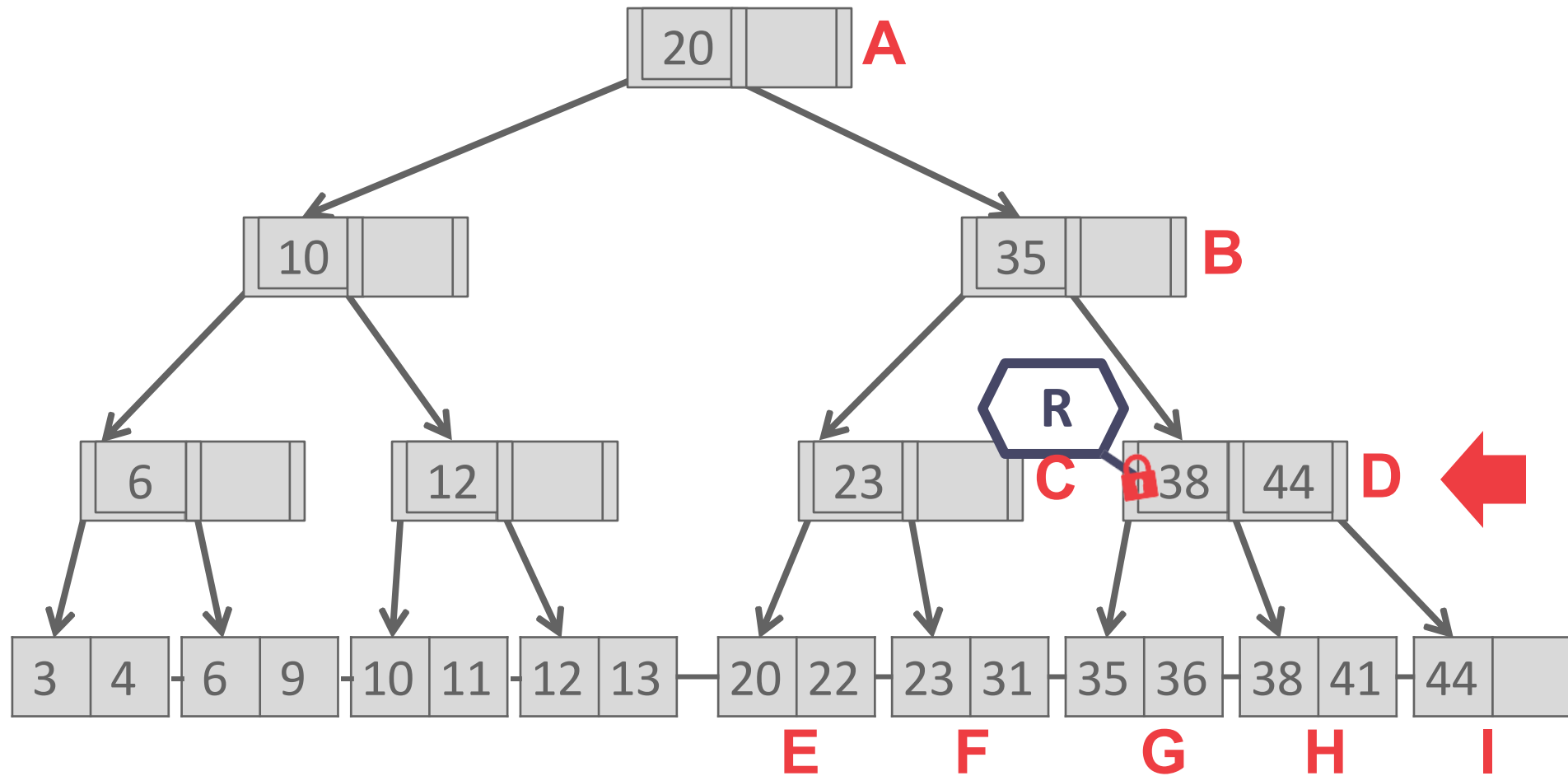
EXAMPLE #2 – DELETE 38



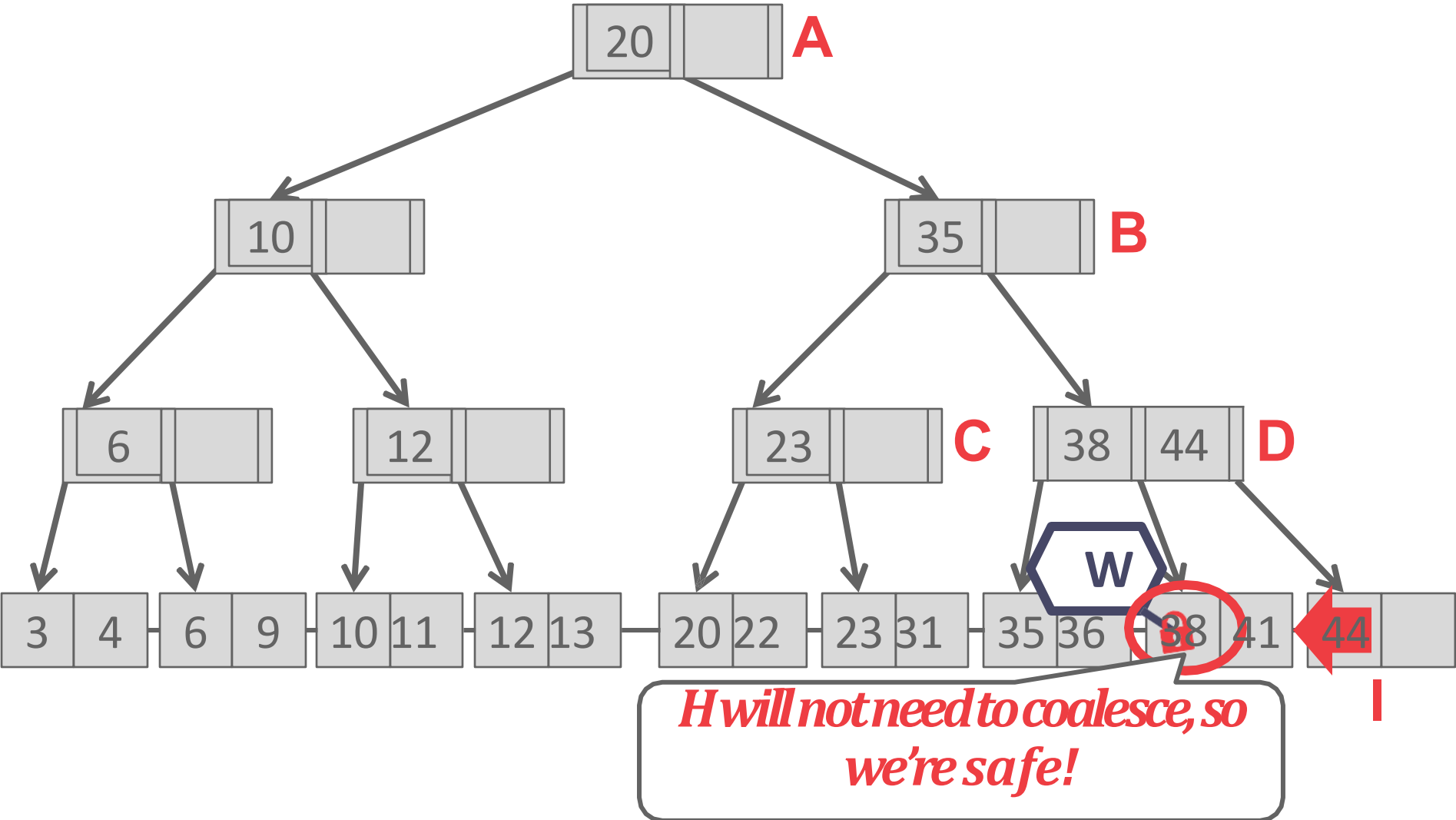
EXAMPLE #2 – DELETE 38



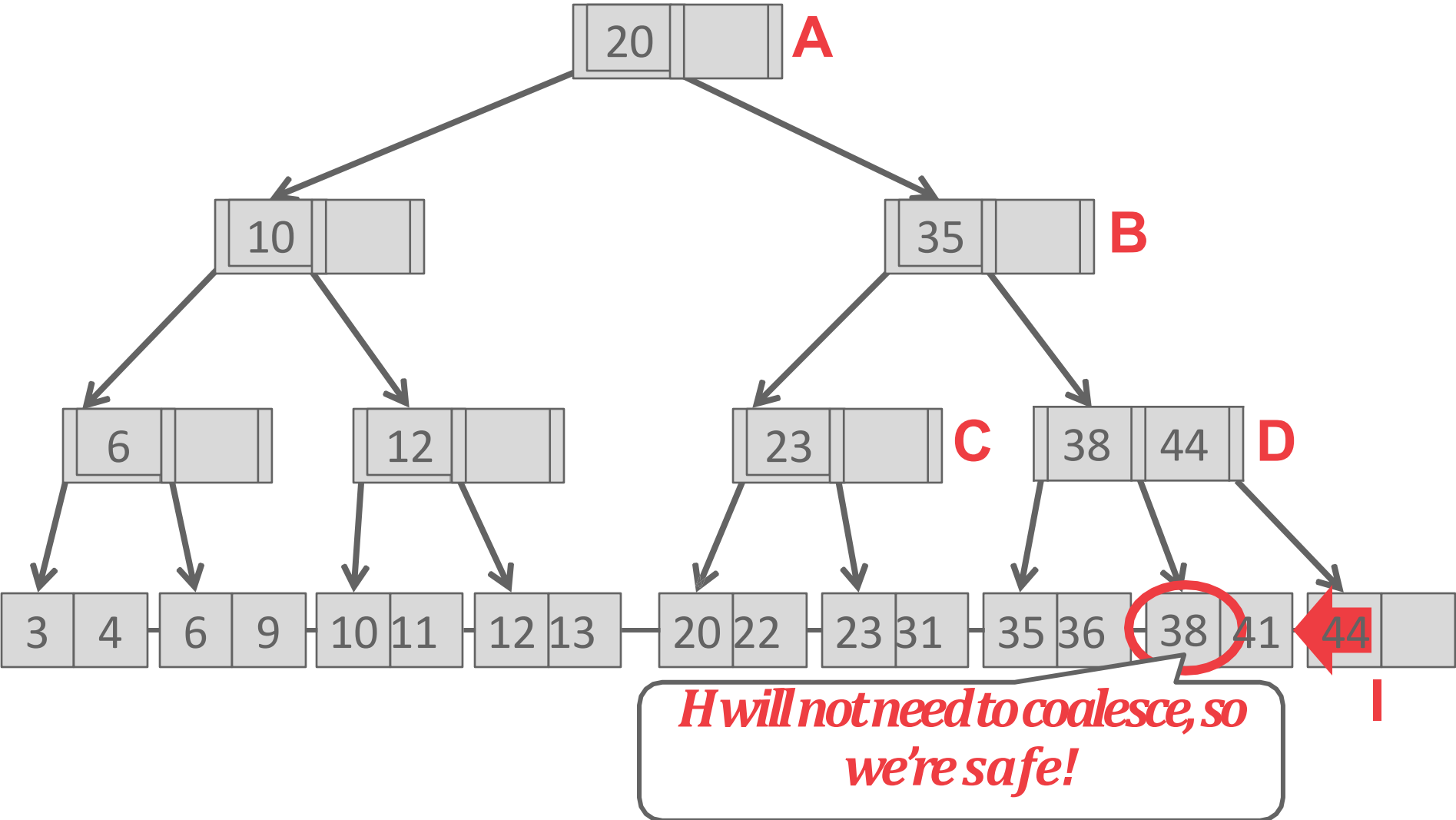
EXAMPLE #2 – DELETE 38



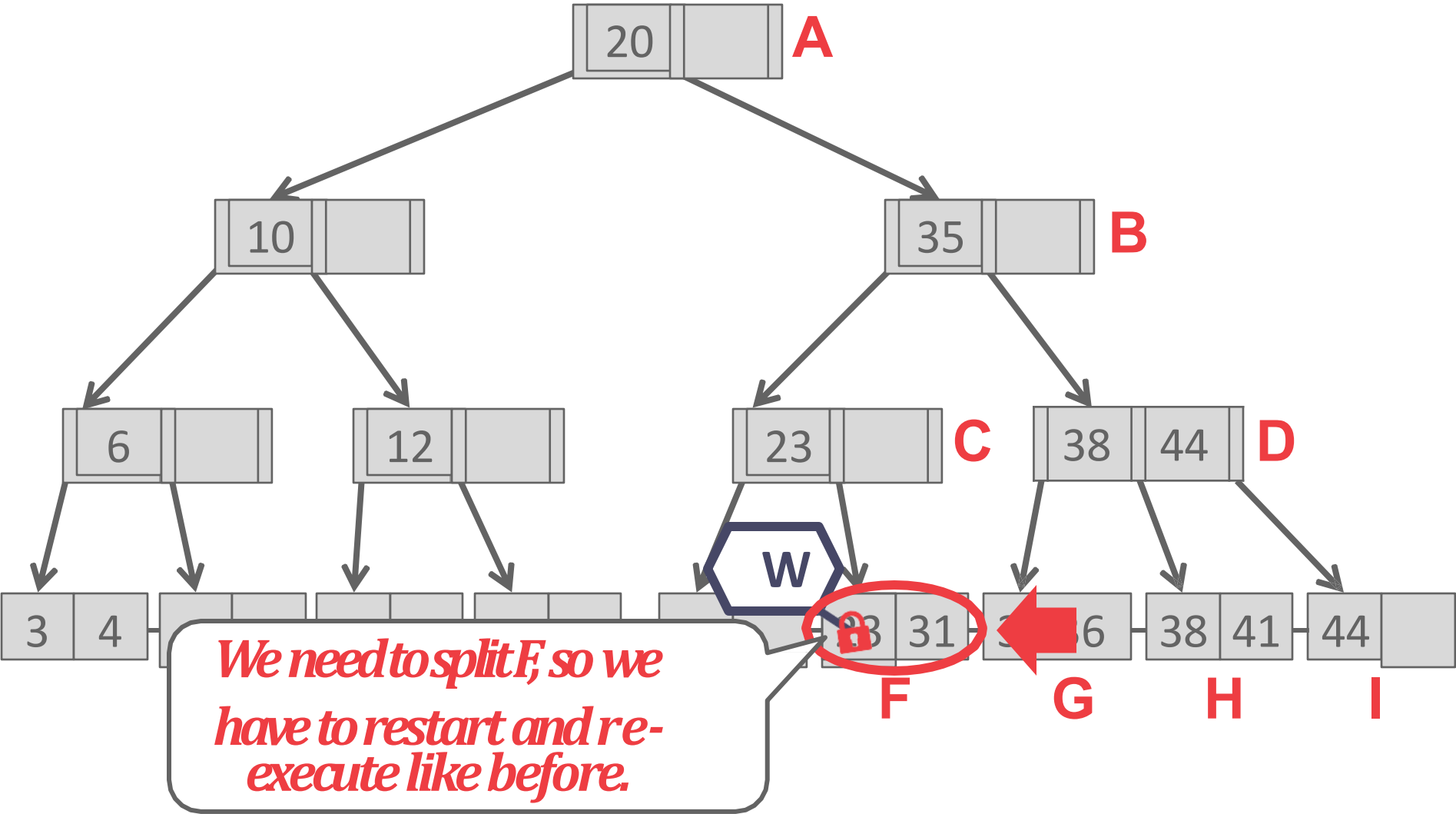
EXAMPLE #2 – DELETE 38



EXAMPLE #2 – DELETE 38



EXAMPLE #4 – INSERT 25



OBSERVATION

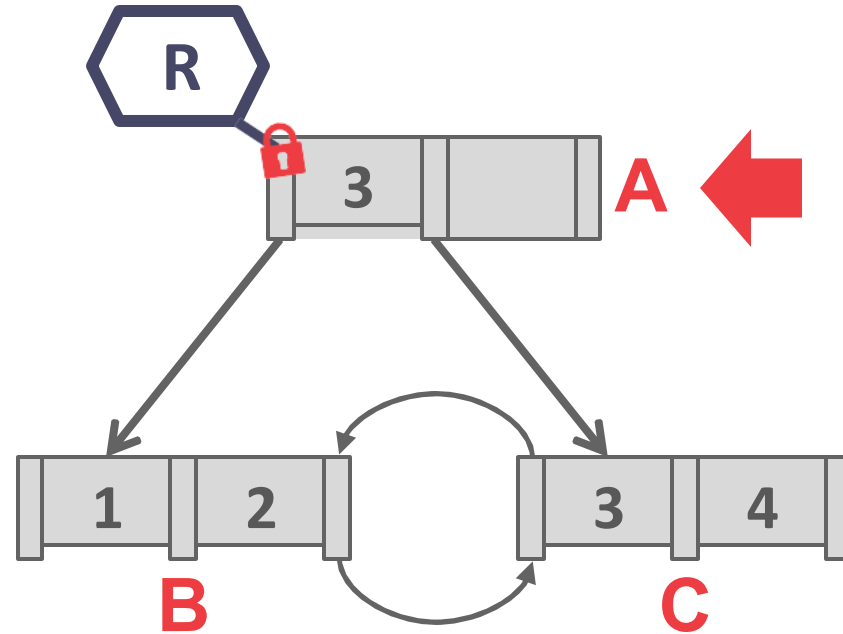
The threads in all the examples so far have acquired latches in a "top-down" manner.

- 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 we want to move from one leaf node to another leaf node?

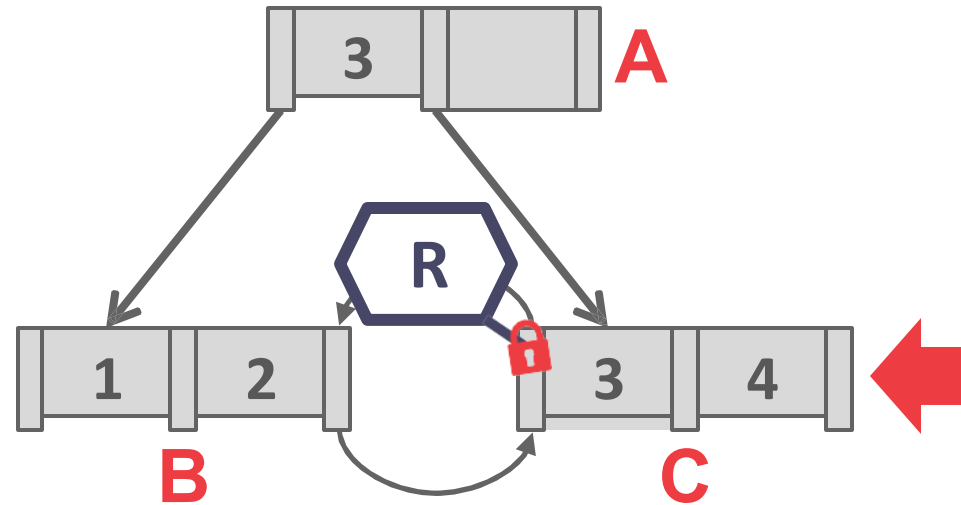
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



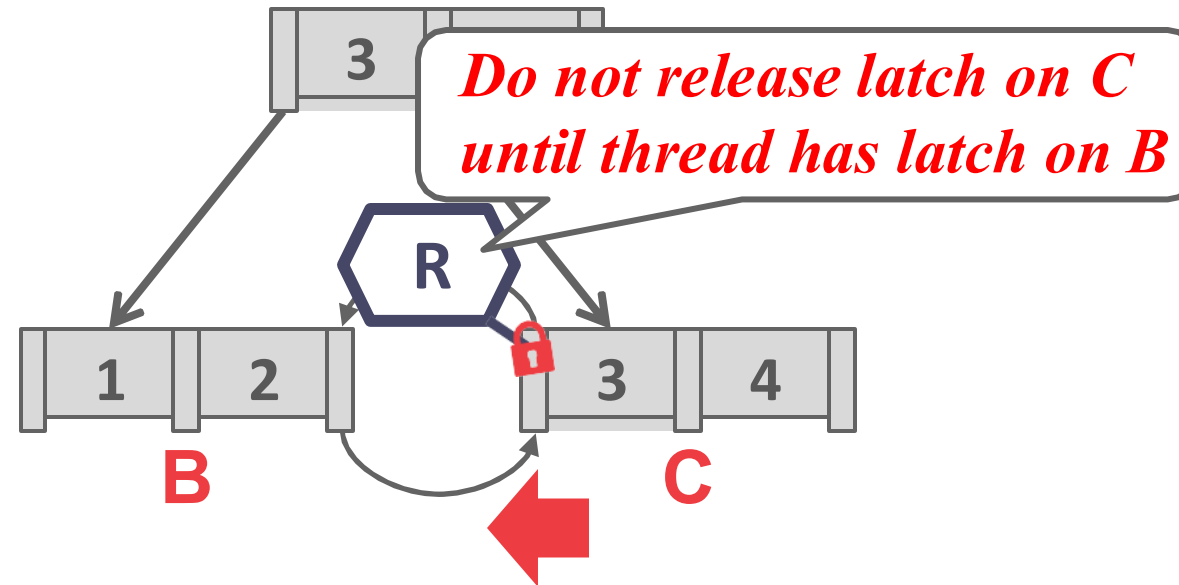
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



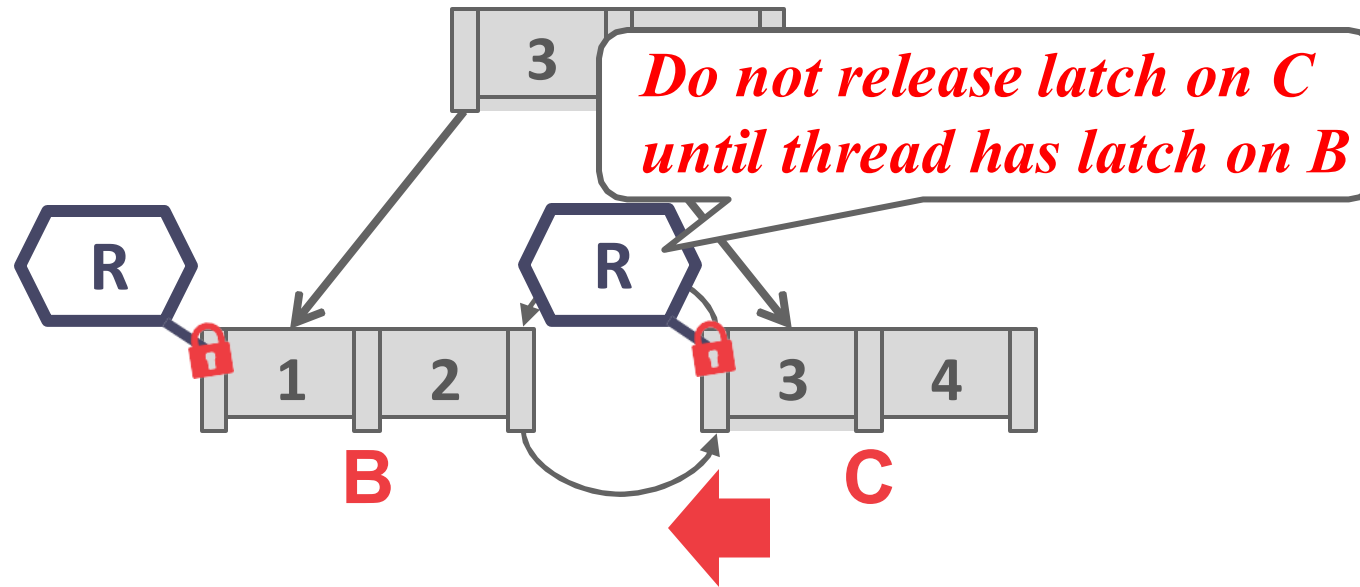
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



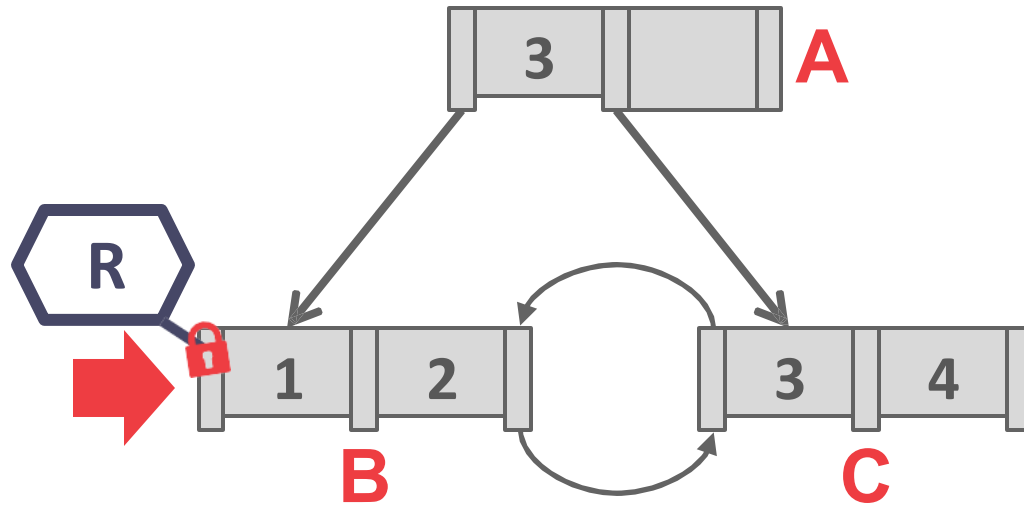
LEAF NODE SCAN EXAMPLE #1

T_1 : Find Keys < 4



LEAF NODE SCAN EXAMPLE #1

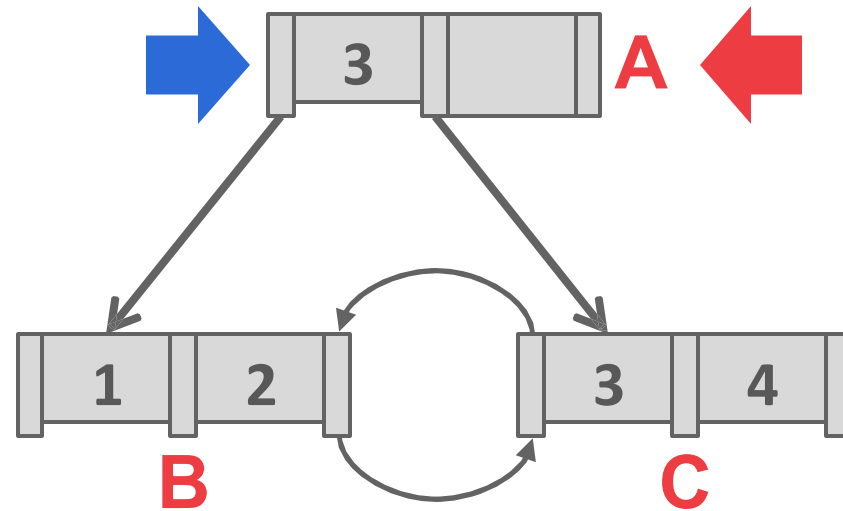
T_1 : Find Keys < 4



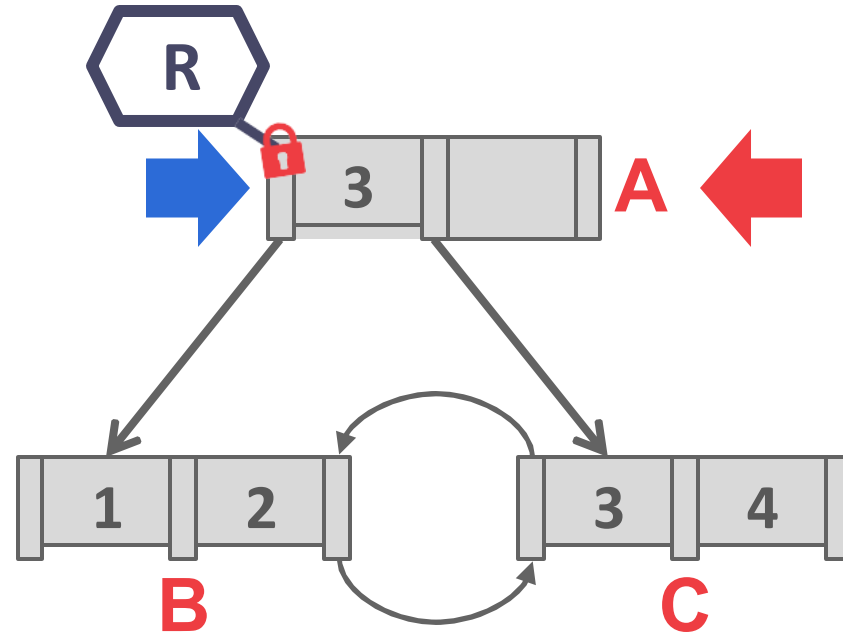
LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2



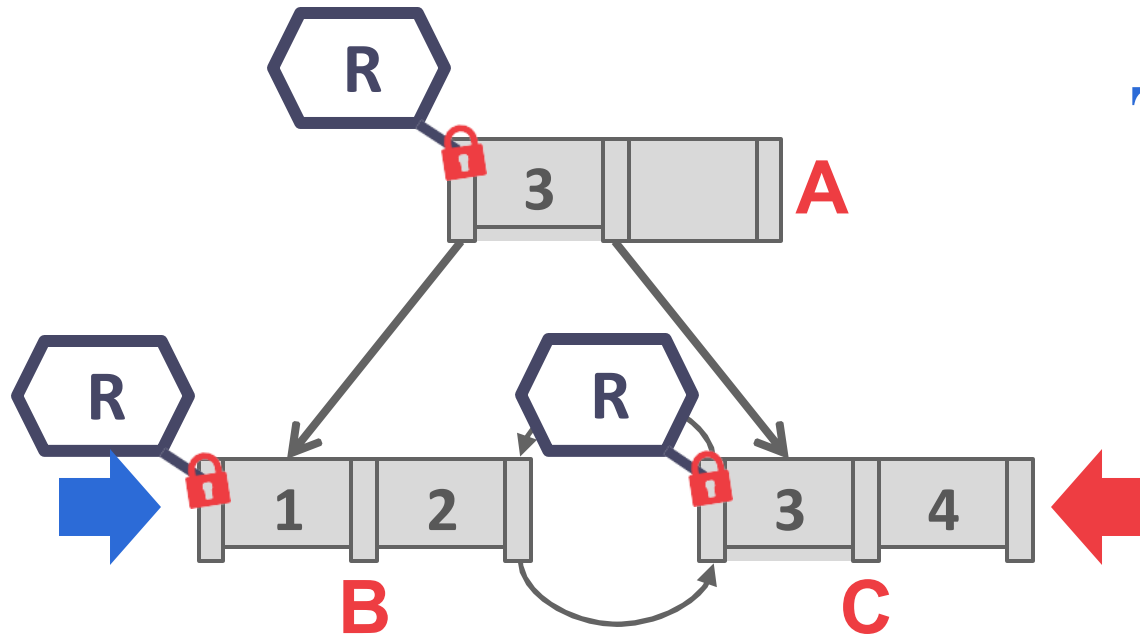
T_1 : Find Keys < 4

T_2 : Find Keys > 1

LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

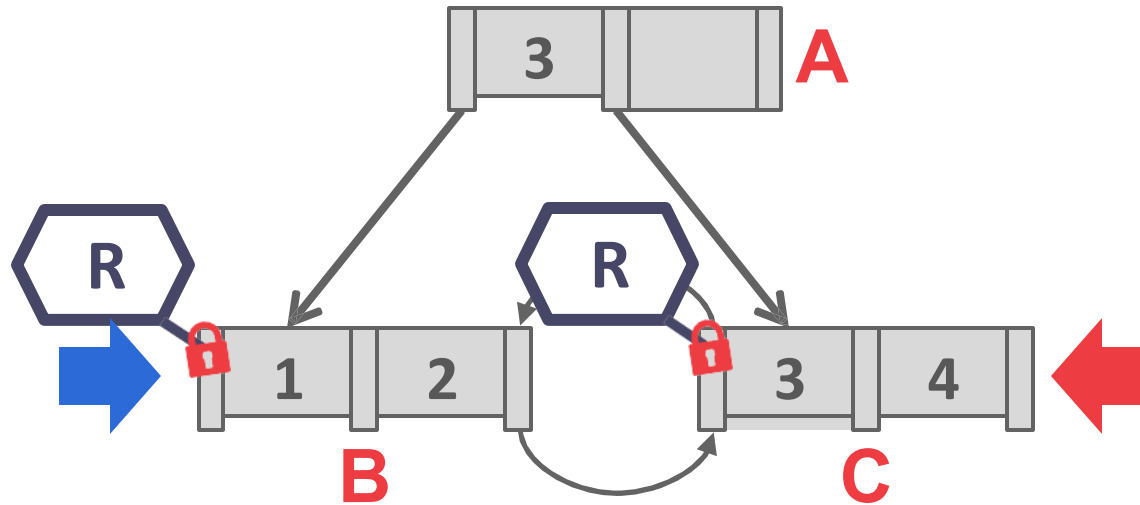
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

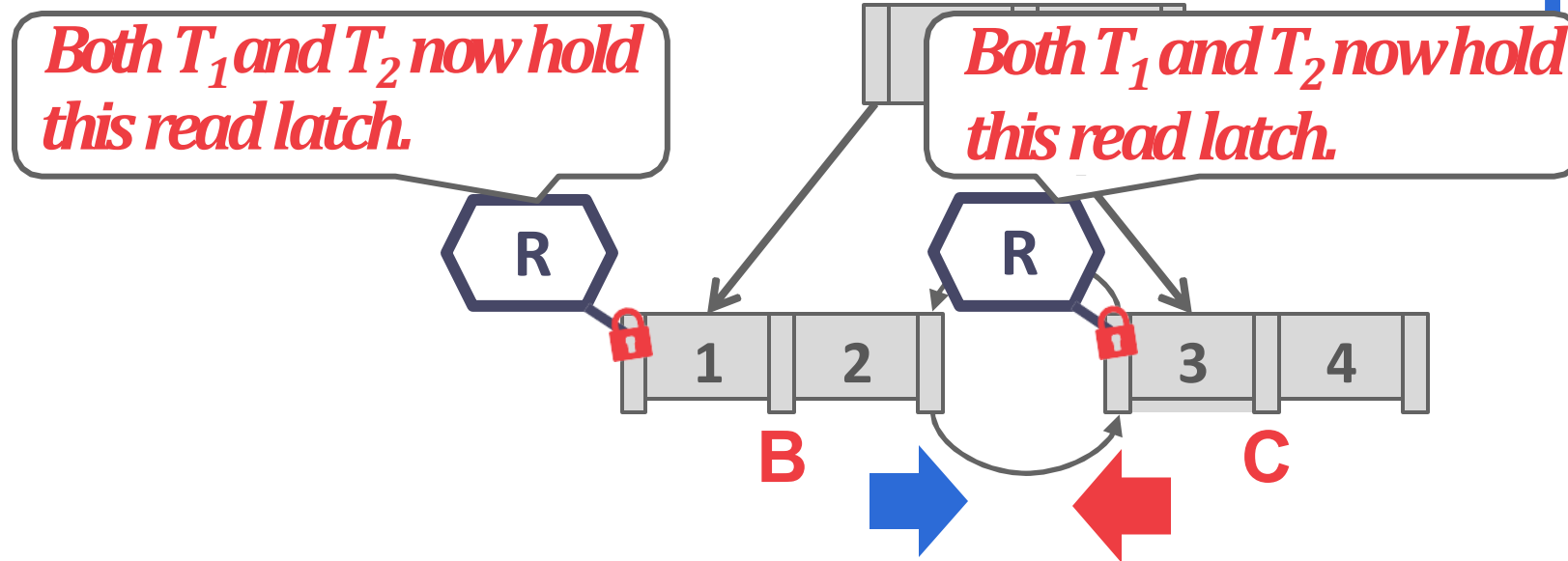
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

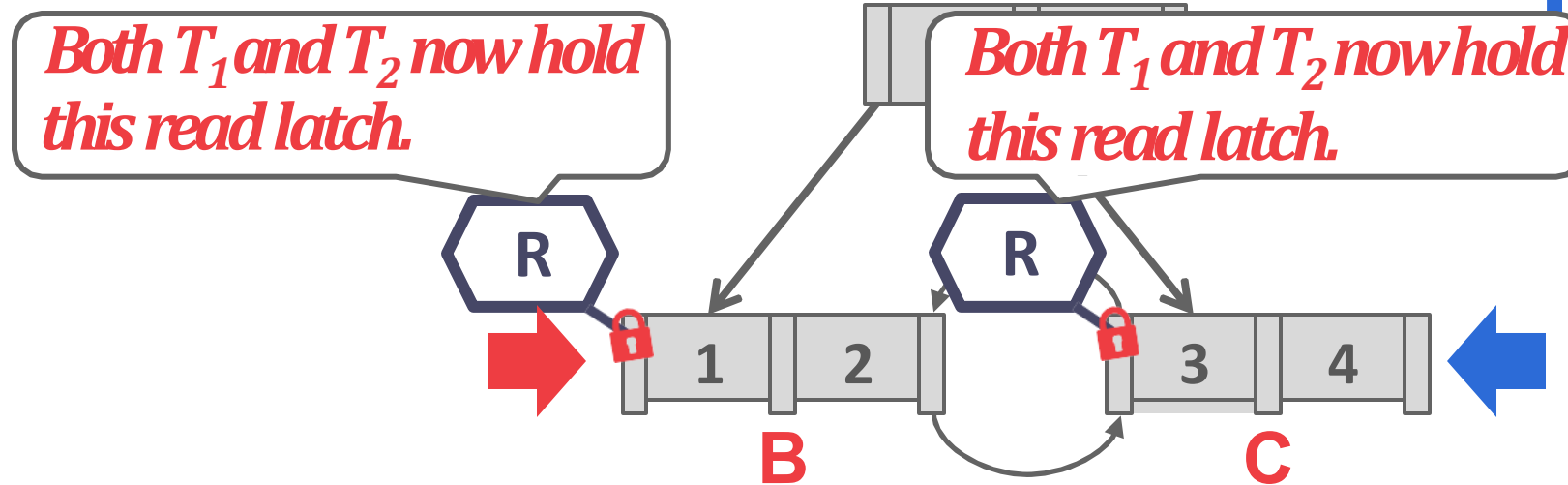
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

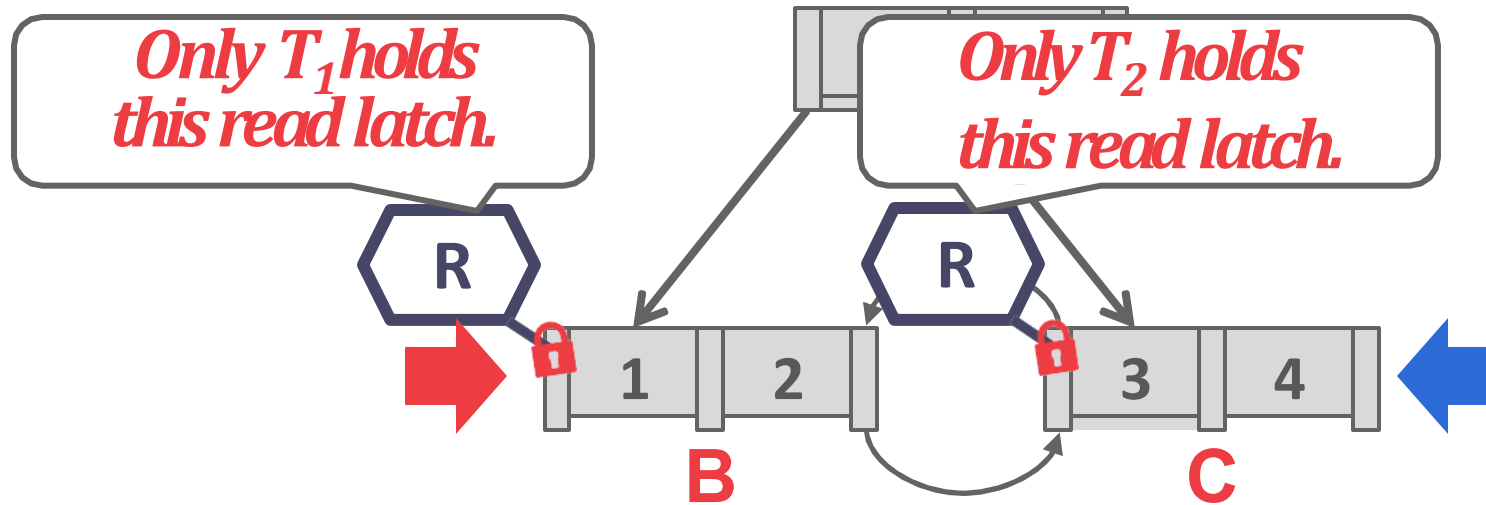
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #2

T_1 : Find Keys < 4

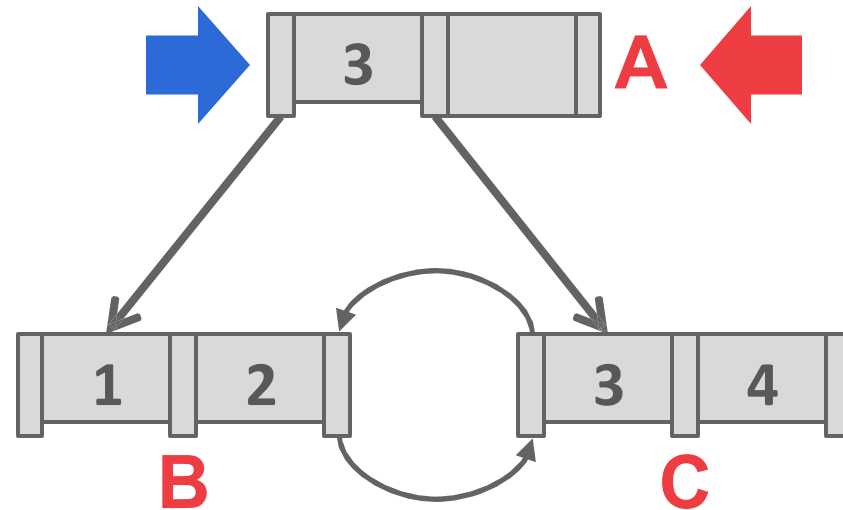
T_2 : Find Keys > 1



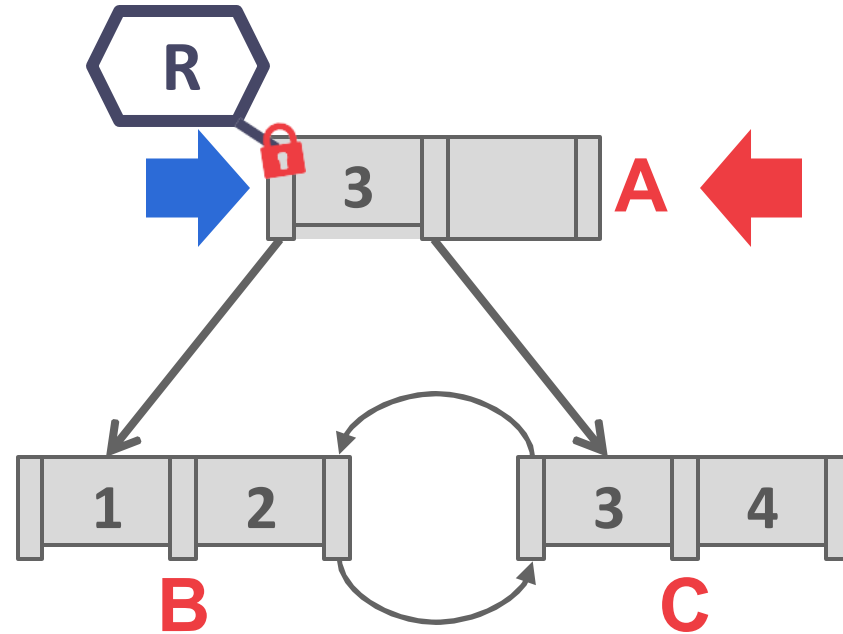
LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3



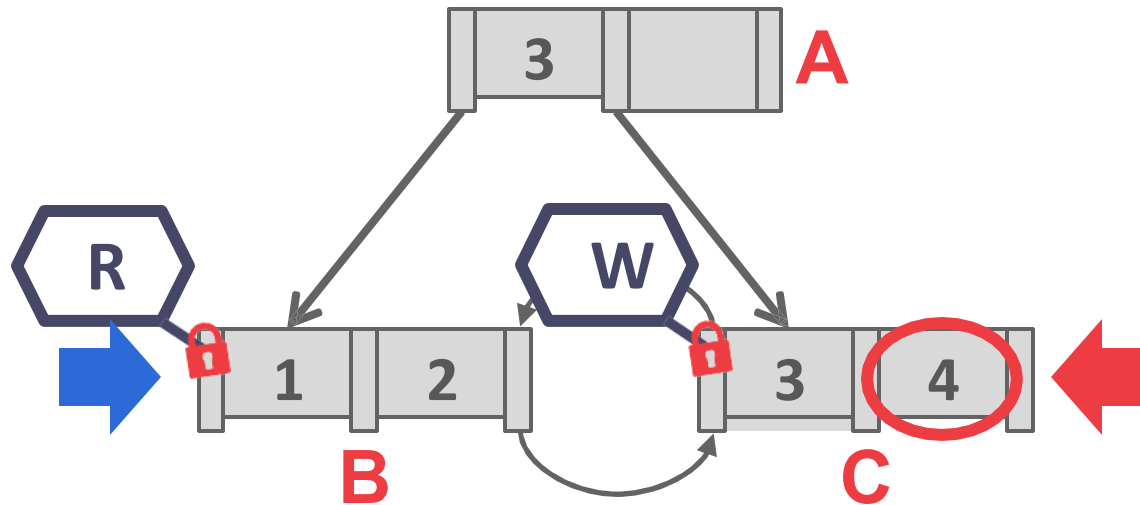
T_1 : Delete 4

T_2 : Find Keys > 1

LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

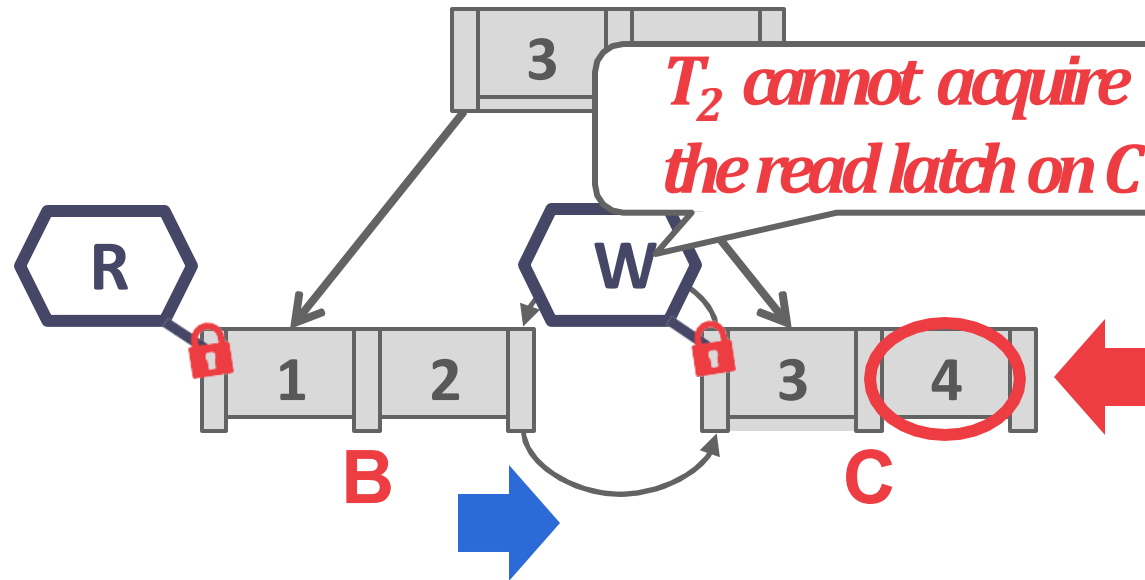
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

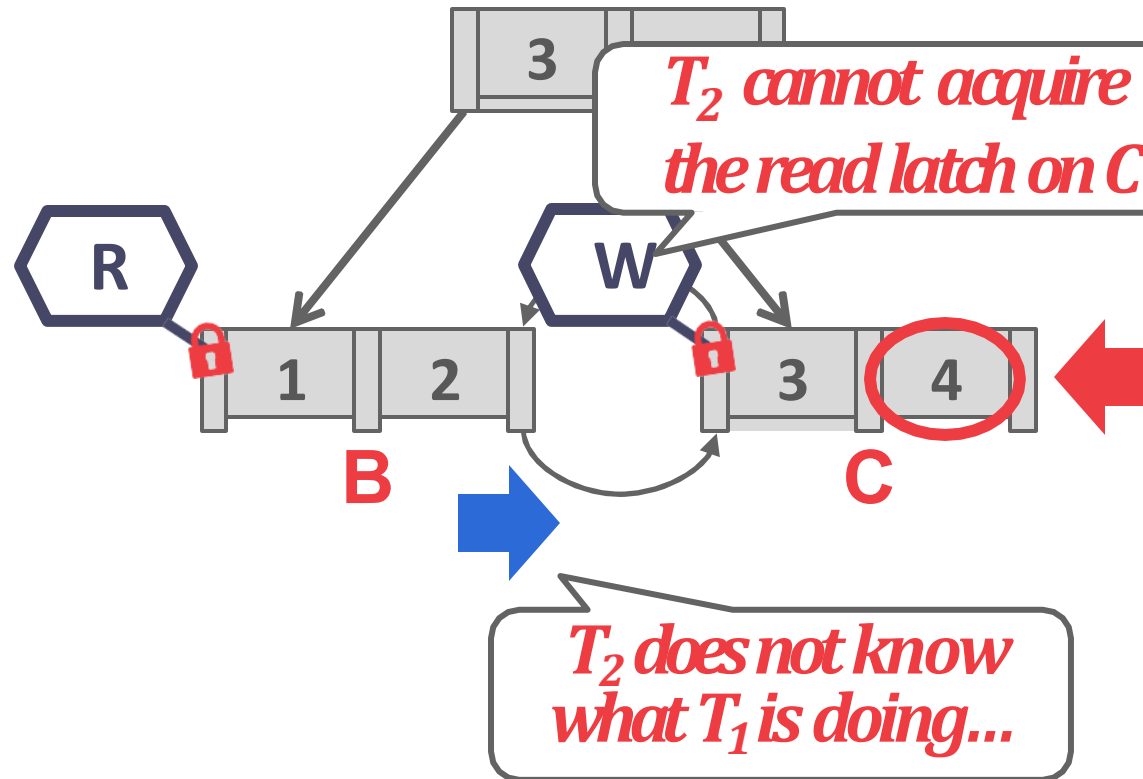
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

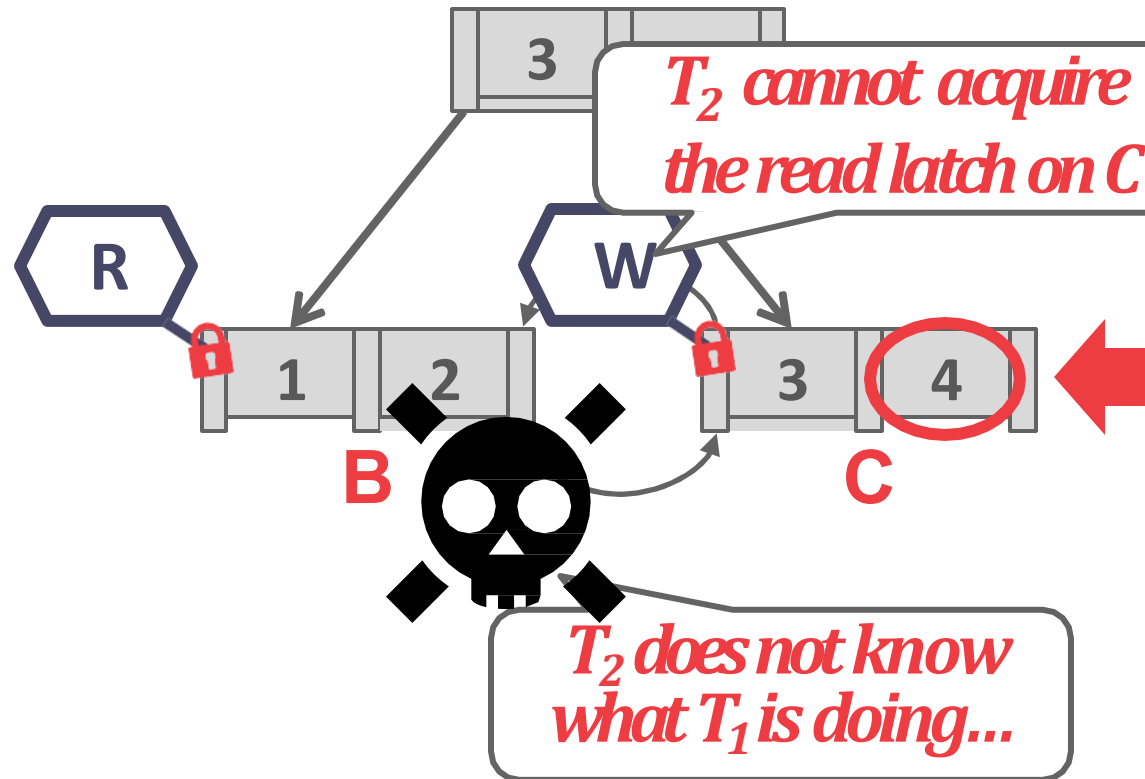
T_2 : Find Keys > 1



LEAF NODE SCAN EXAMPLE #3

T_1 : Delete 4

T_2 : Find Keys > 1



LEAF NODE 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.

DELAYED PARENT UPDATES

Every time a leaf node overflows, we must update at least three nodes.

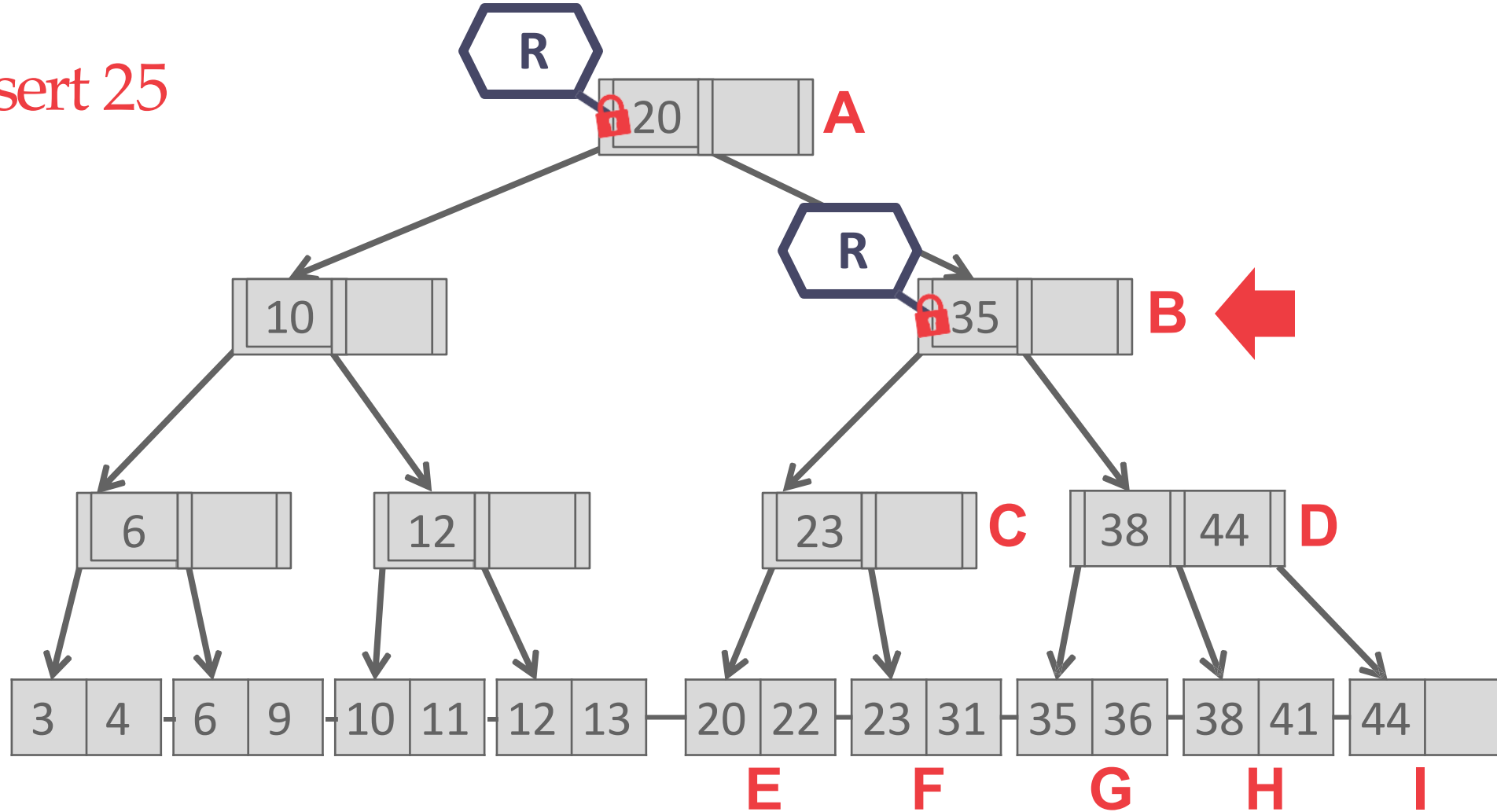
- The leaf node being split.
- The new leaf node being created.
- The parent node.

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

EXAMPLE #4

- INSERT 25

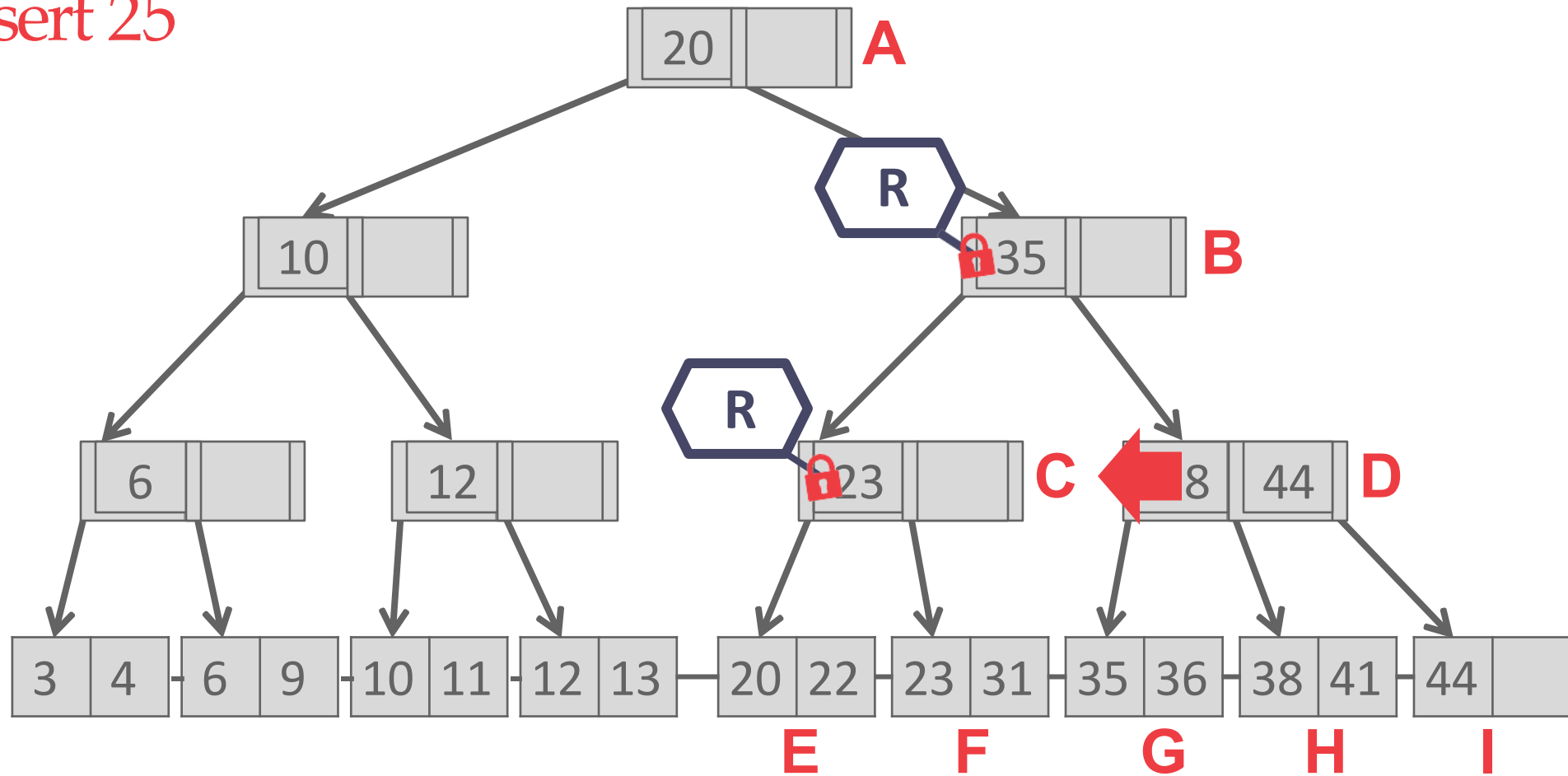
T_1 : Insert 25



EXAMPLE #4

– INSERT 25

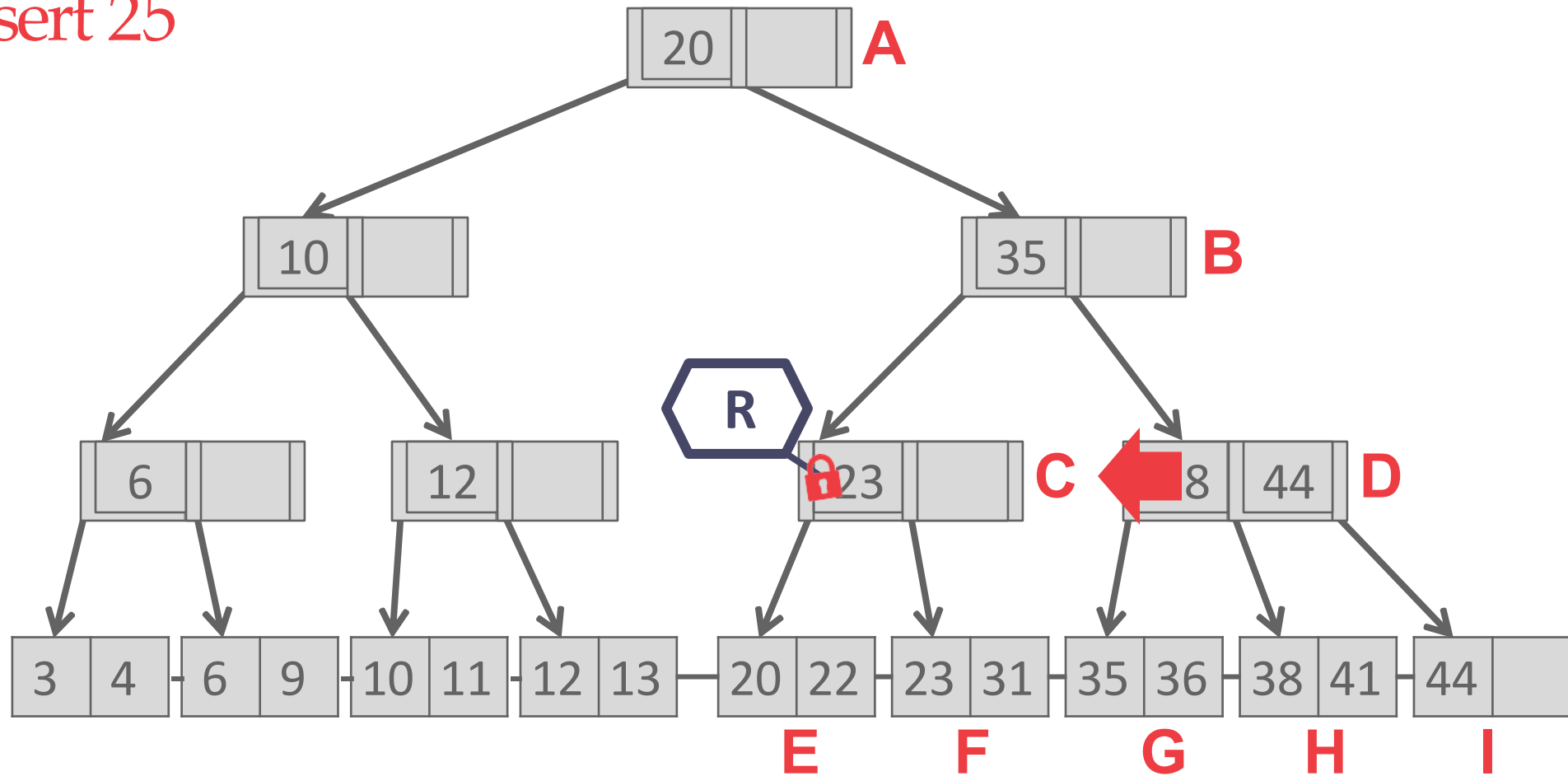
T_1 : Insert 25



EXAMPLE #4

– INSERT 25

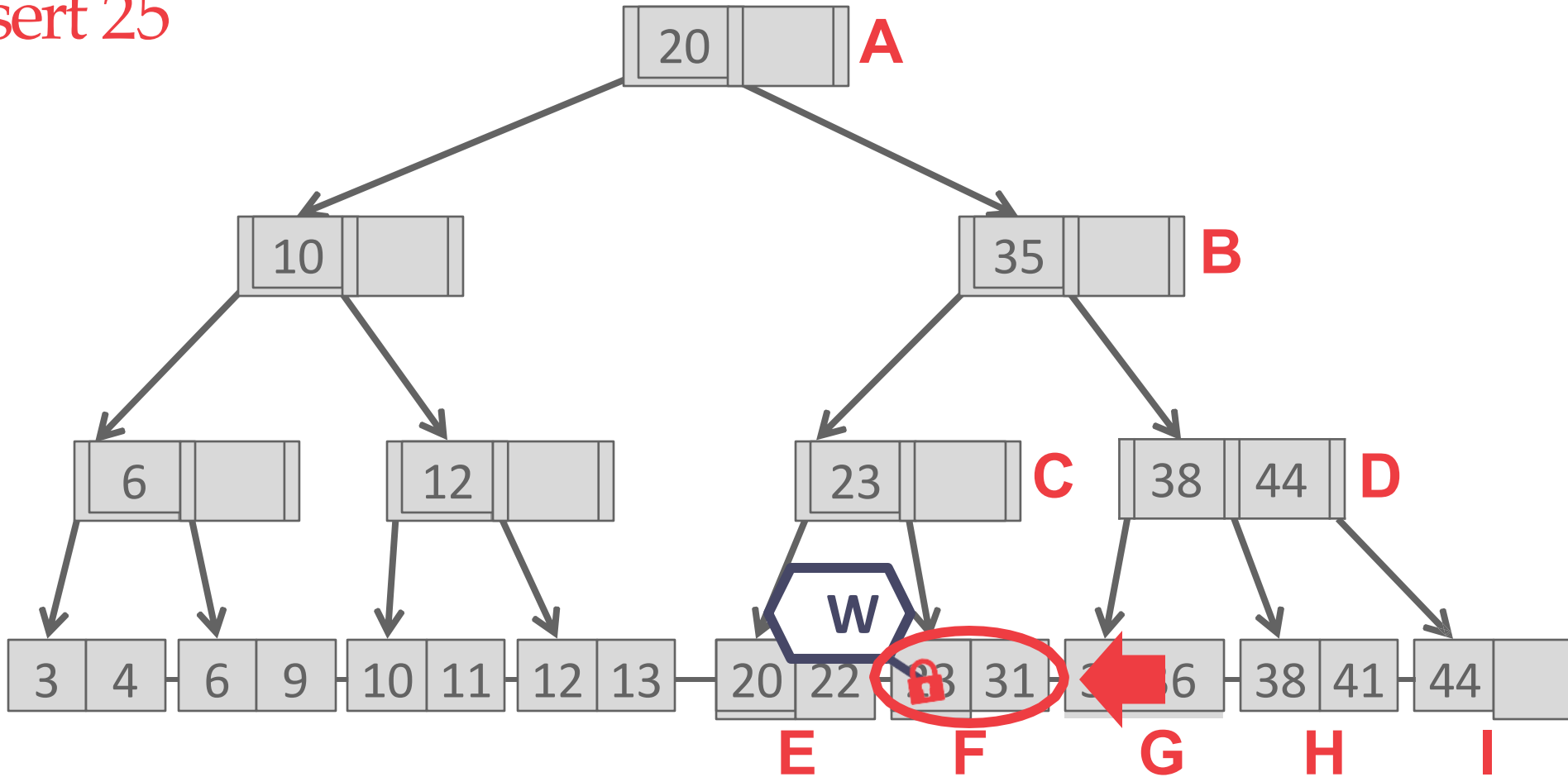
T_1 : Insert 25



EXAMPLE #4

– INSERT 25

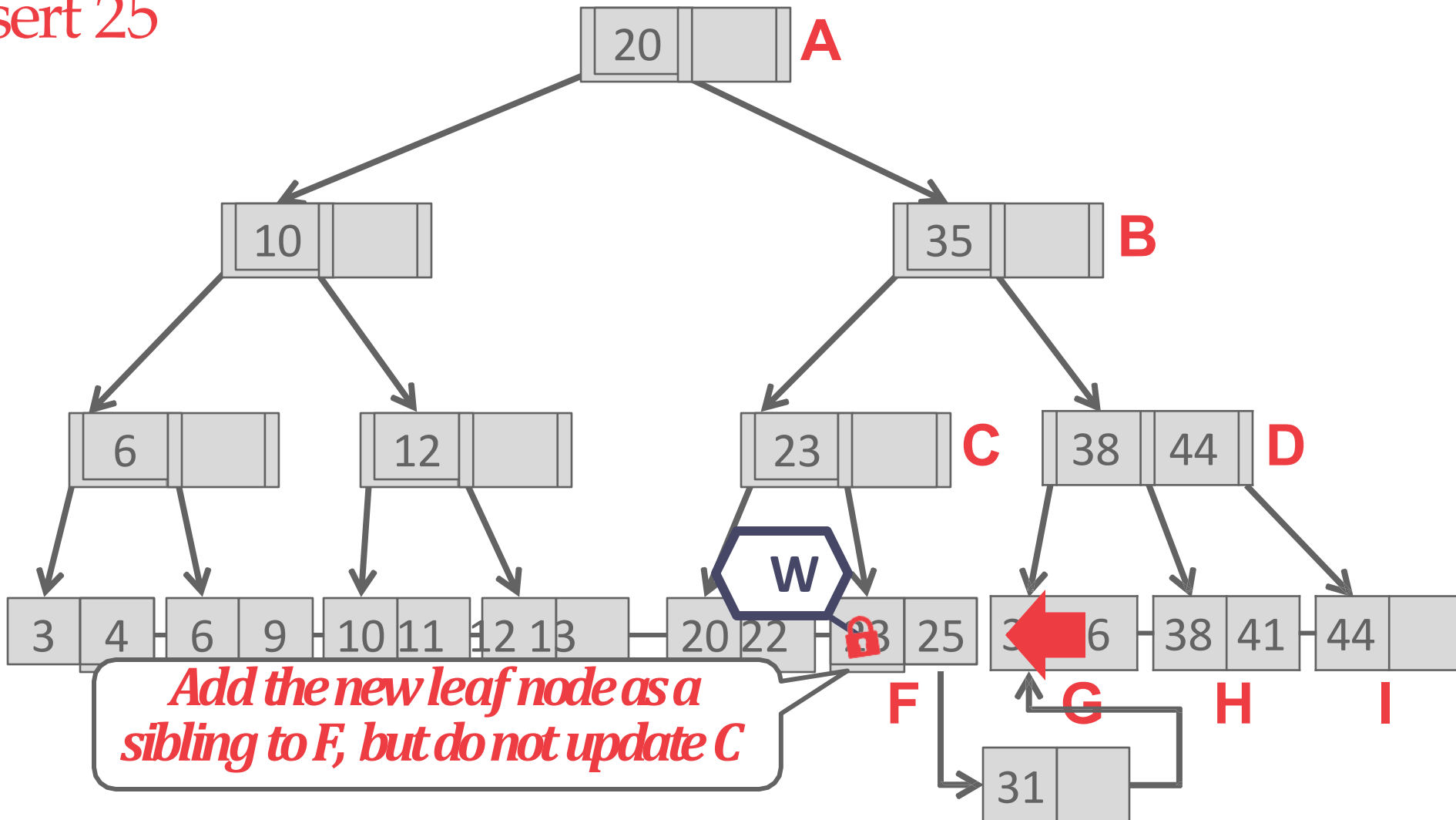
T_1 : Insert 25



EXAMPLE #4

– INSERT 25

T_1 : Insert 25

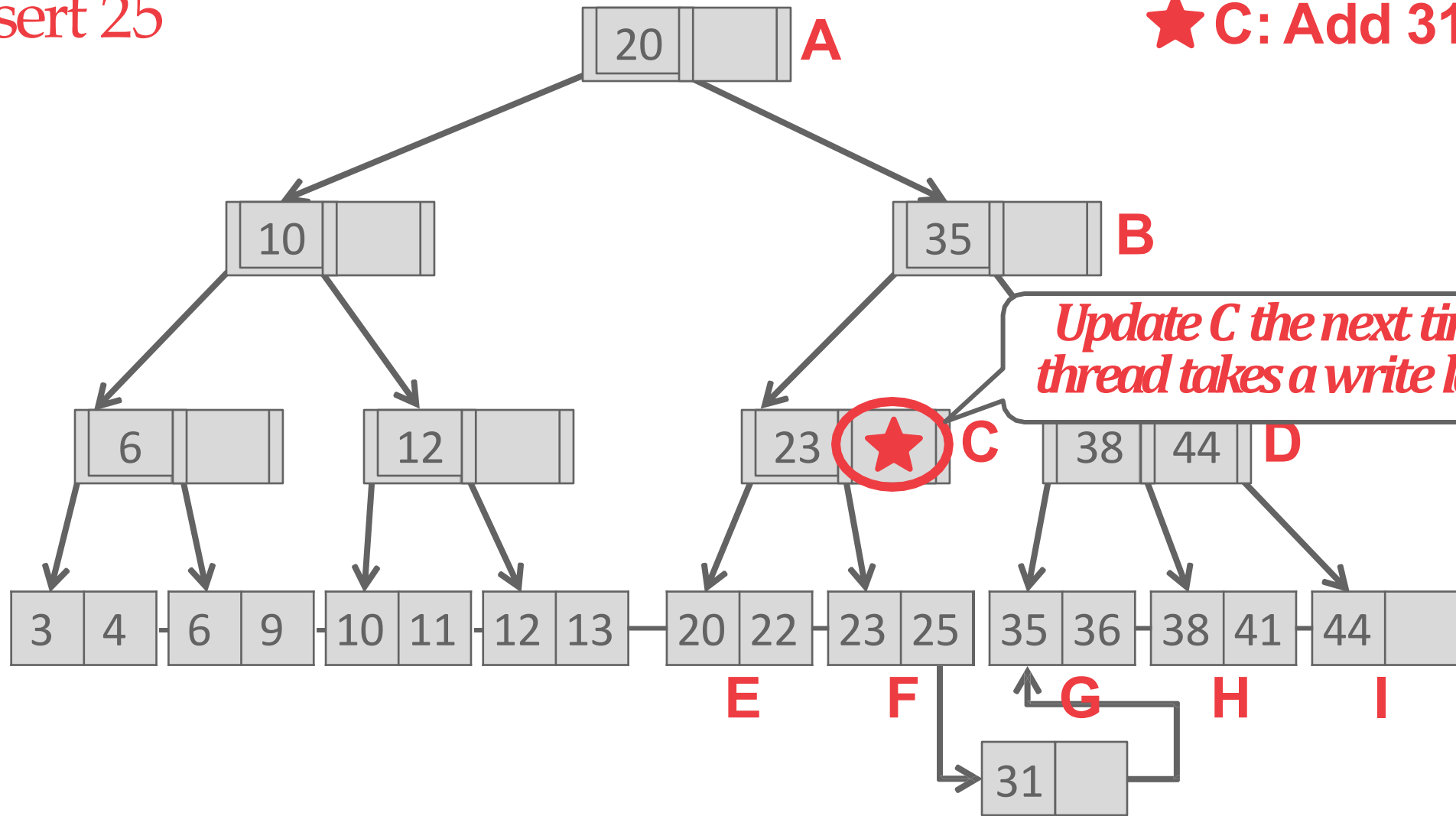


EXAMPLE #4

– INSERT 25

T_1 : Insert 25

★ C: Add 31

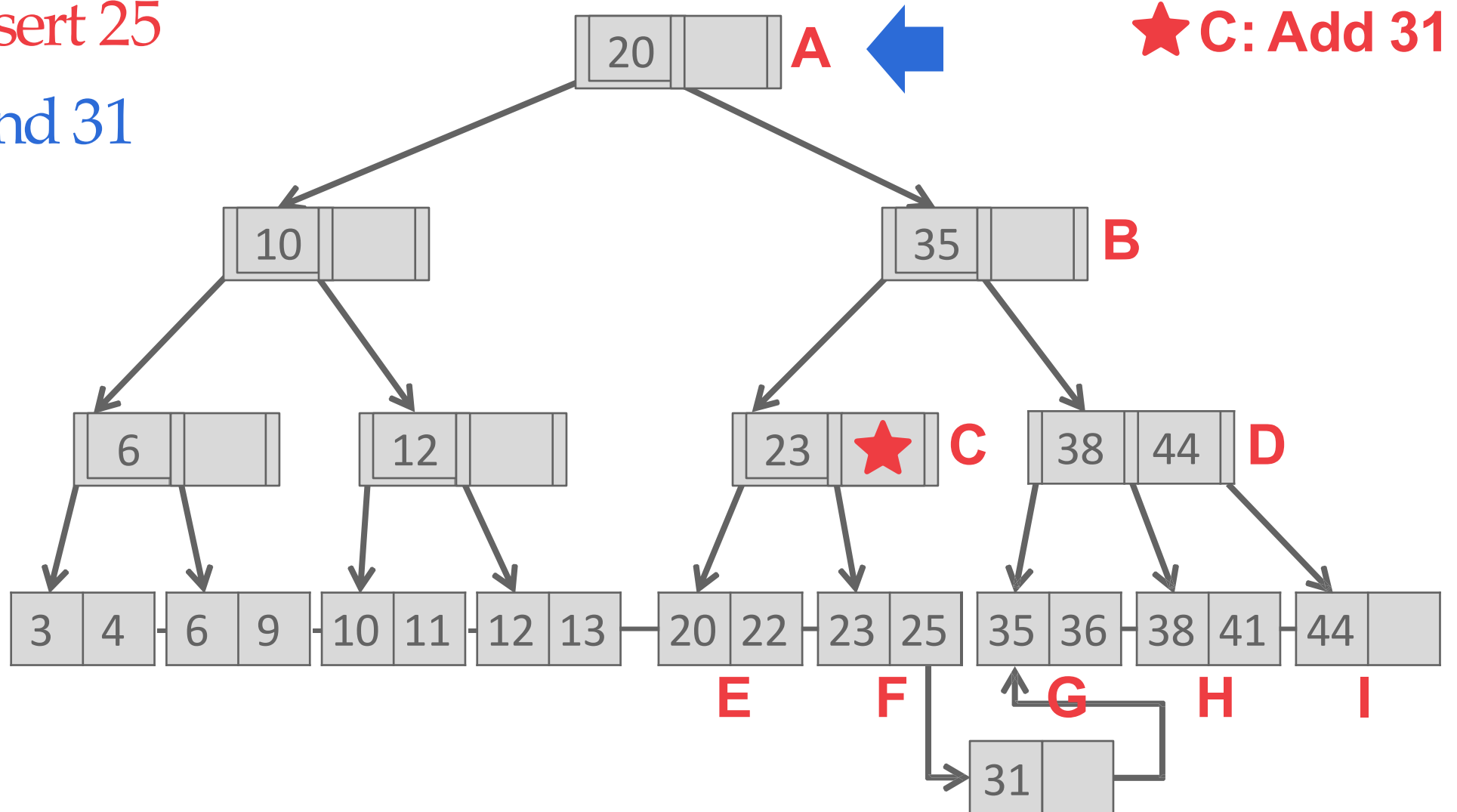


EXAMPLE #4

– INSERT 25

T_1 : Insert 25

T_2 : Find 31



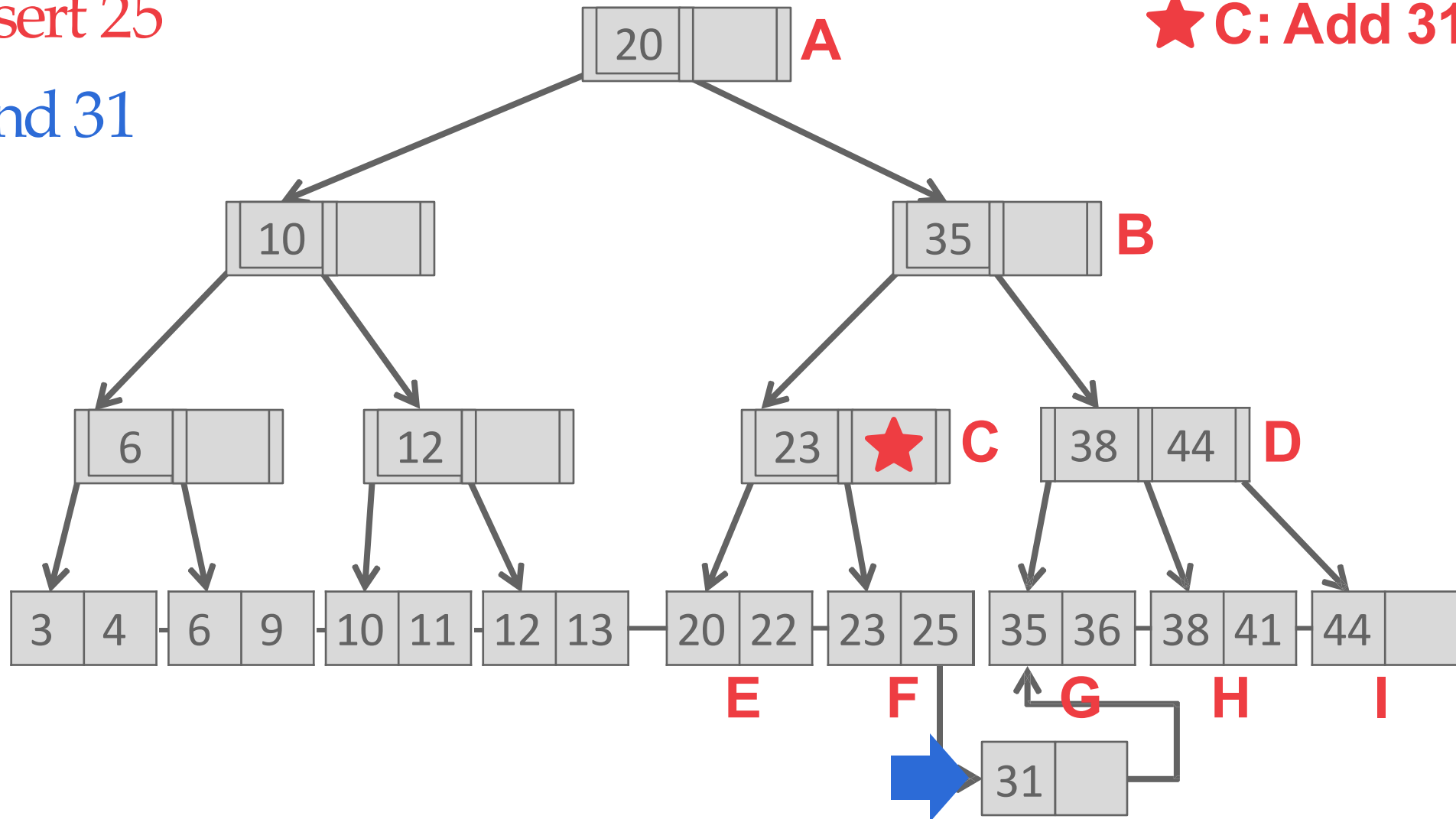
EXAMPLE #4

– INSERT 25

T_1 : Insert 25

T_2 : Find 31

★ C: Add 31



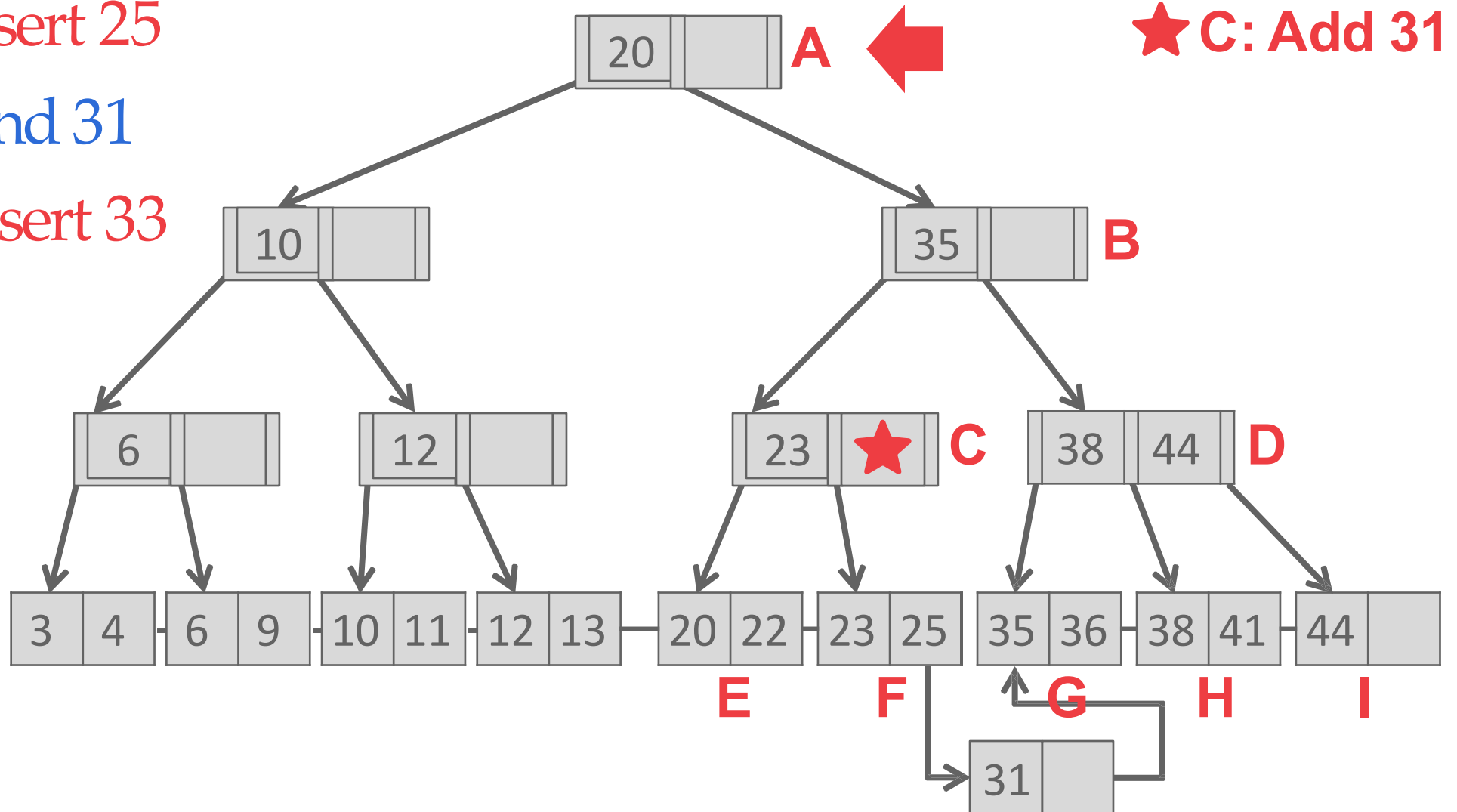
EXAMPLE #4

– INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33



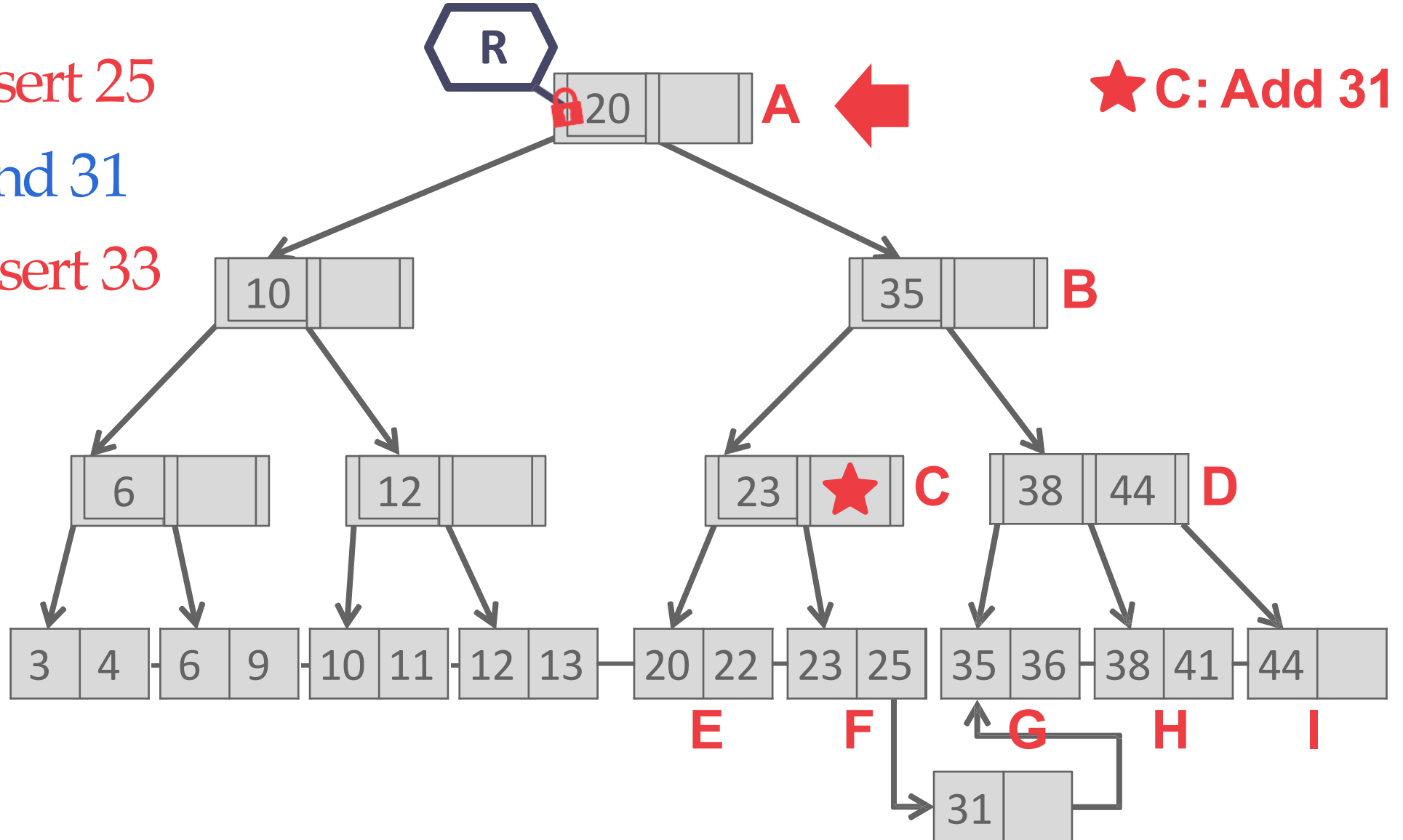
EXAMPLE #4

- INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33



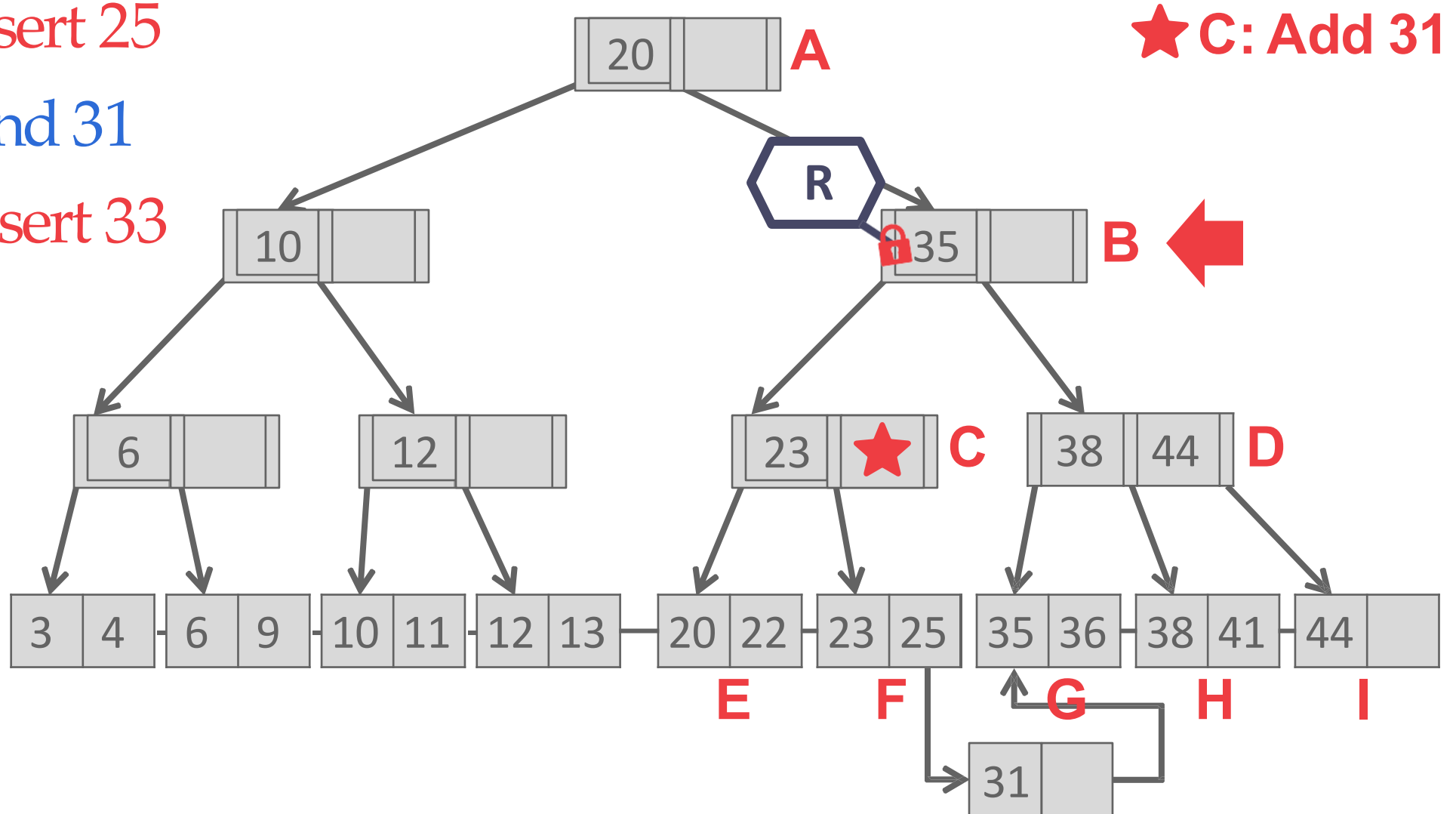
EXAMPLE #4

– INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33



EXAMPLE #4

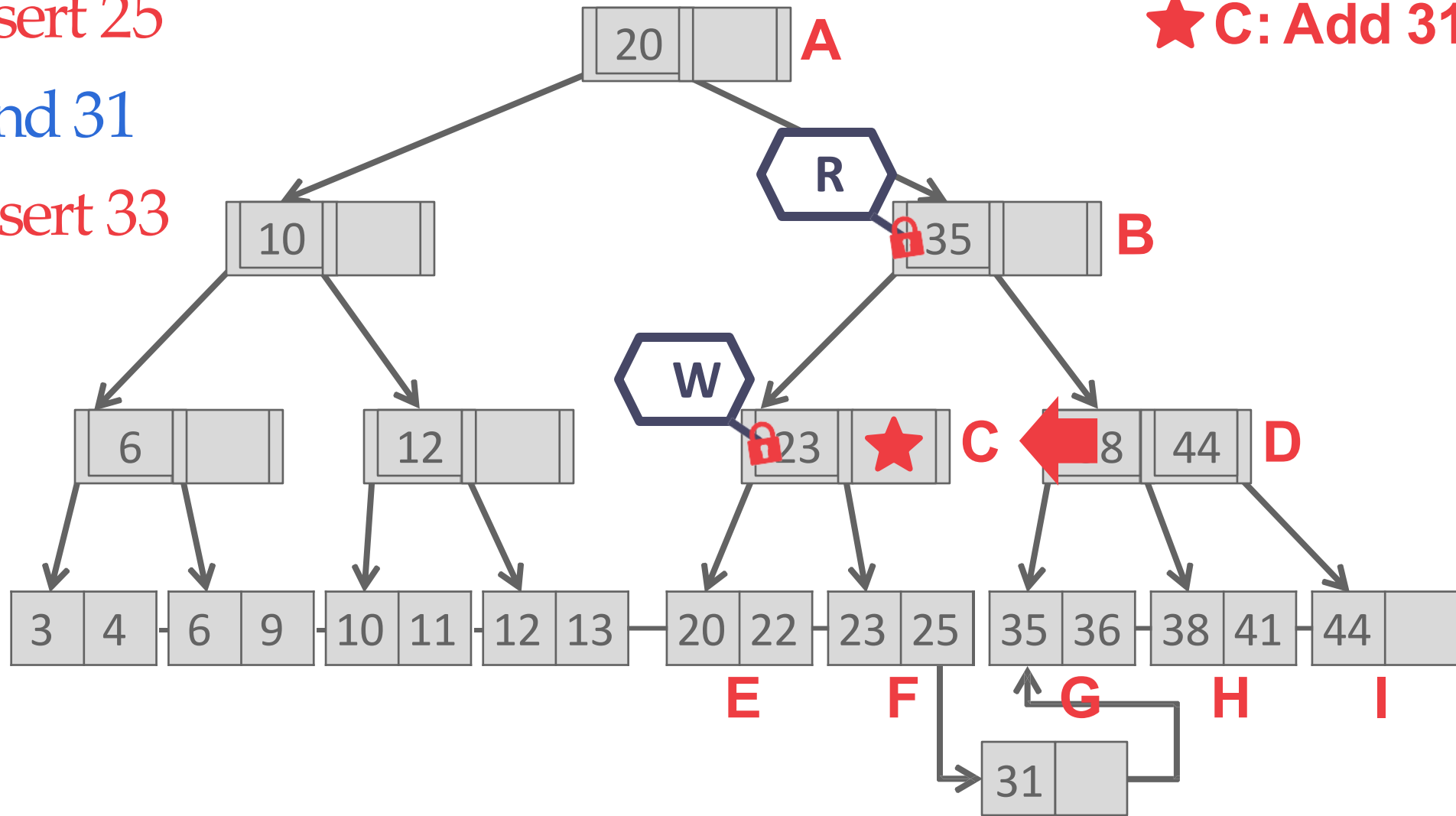
– INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33

★ C: Add 31



EXAMPLE #4

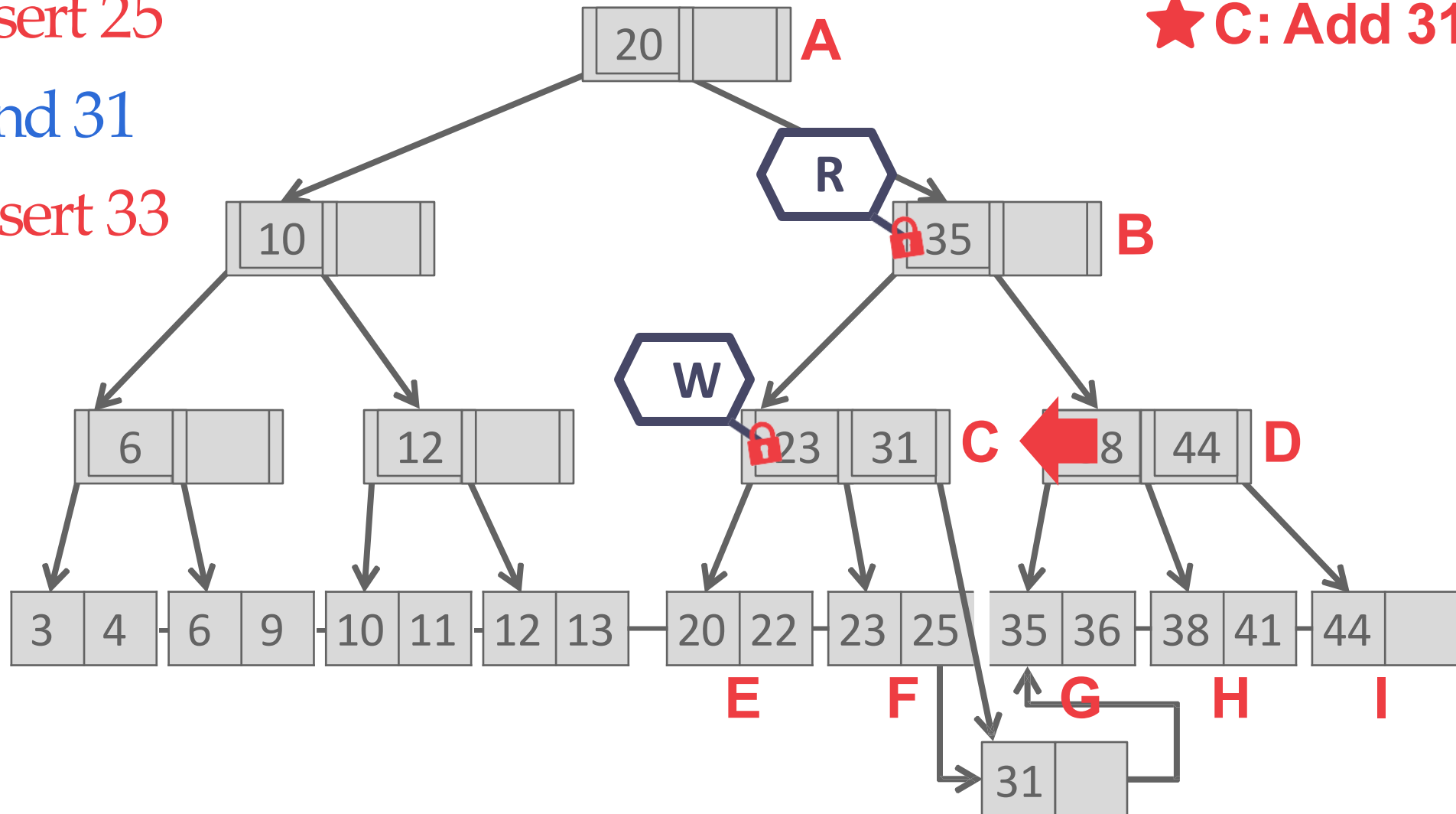
– INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33

★ C: Add 31



EXAMPLE #4

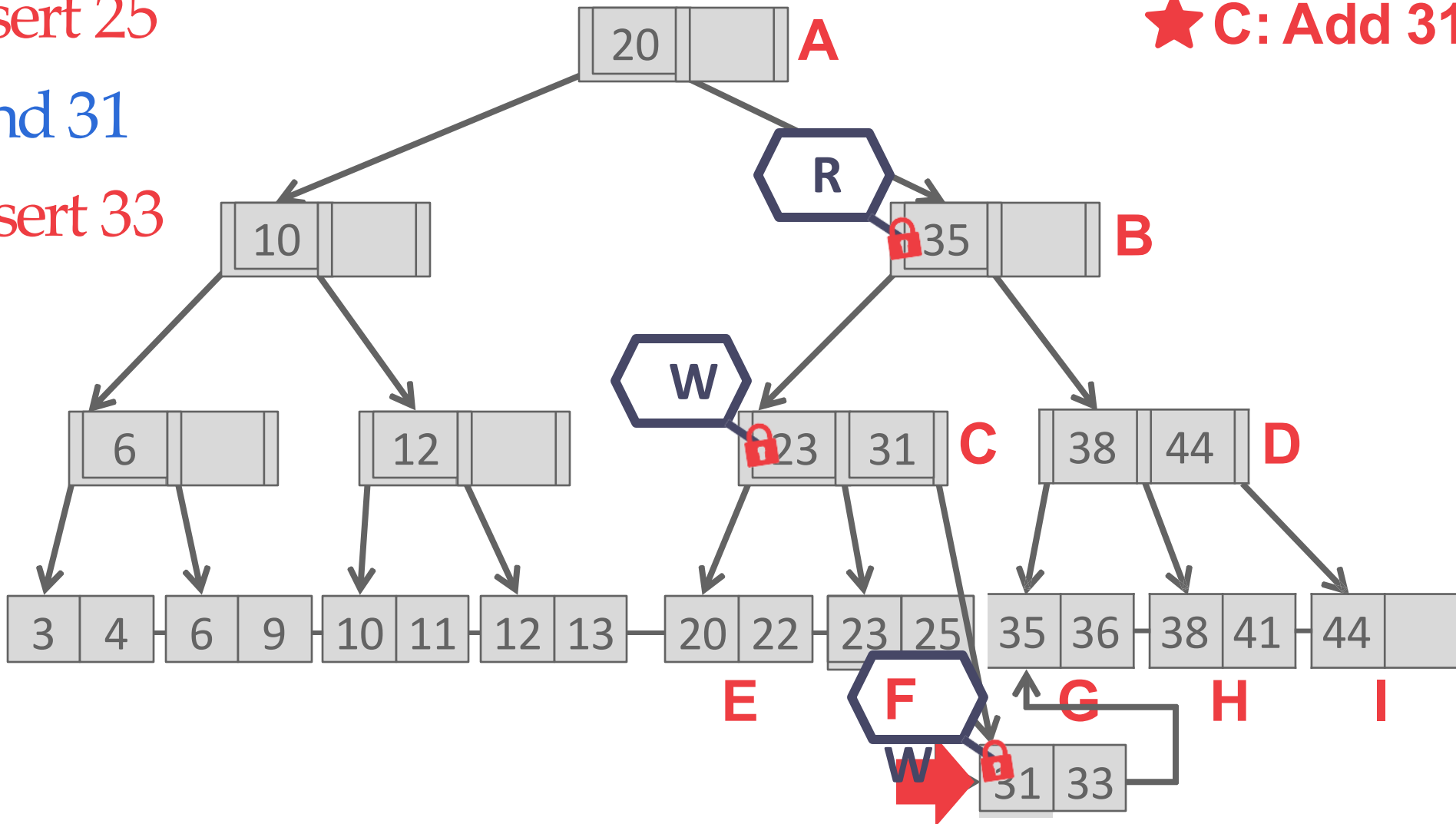
– INSERT 25

T_1 : Insert 25

T_2 : Find 31

T_3 : Insert 33

★ C: Add 31



Next lecture

- Locking and transactions
- MVCC

Potential course projects: Concurrency

- Faster dynamic graph processing
 - Terrace with parallel PMA/BP-tree
- Implement concurrency in on-disk B+-tree
- Implement transactions in a key-value store

Potential course projects: Filters

- Adaptive filter-based joins
- Adaptive filter-based minimal perfect hashing
- Summary cache redesign using maplet. (MemcacheD/Redis)

Potential course projects: Indexing

- In-memory write-optimization

Potential course projects: Query Optimization

- DuckDB query optimization
- Postgres query optimization - processing

Potential course projects: GPU

- GPU-CPU caching system
- GPU-CPU write optimization

Potential course projects: Logging/Recovery

- Implement logging and recovery in a key-value store