

CS 6400 A

# Database Systems Concepts and Design

---

Lecture 11  
09/29/25

# Announcements

## Midterm grade released

- After adjustment: mean: 80.3, median: 83.5, std: 14.9
- Answer key on Canvas; regrade request open till next Monday (Oct 6)

## Assignment 2 will be released tonight

- Programming assignment on in-memory data layout; due **Oct 20**
- Please START EARLY on this assignment! Autograder takes ~15min to run

Project proposal due this Wednesday (Oct 1). No late days.

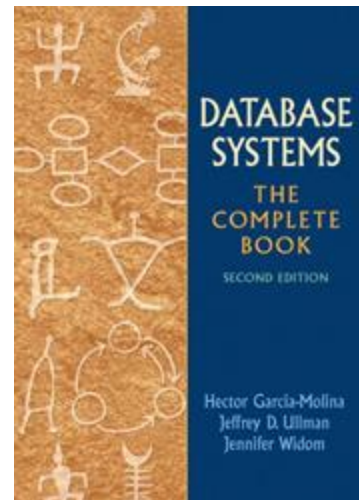
# Agenda

1. B+-Tree Basics
2. B+-Tree Operations
3. Cost Model

# Reading Materials

Database Systems: The Complete Book (2nd edition)

- Chapter 14.2: B-Tree



Acknowledgement: The following slides have been adapted from EE477 (Database and Big Data Systems) taught by Steven Whang, CS145 (Intro to Big Data Systems) taught by Peter Bailis, and CS 6530 (Advanced Database Systems) taught by Prashant Pandey.

# 1. B+-Tree Basics

# B Tree/B+ Tree Overview

They are search trees

- Common misunderstanding: B does not mean binary!
- More general index structure that is commonly used in commercial DBMS's

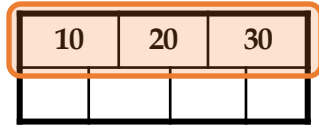
Idea in B Trees:

- Balanced, height adjusted tree
- Stores data (keys and values) in all nodes (both internal and leaf)
- Leaf nodes are independent; no connections between them

Idea in B+ Trees:

- Stores data only in leaf nodes; make leaves into a linked list (for range queries)
- Most popular variant (our focus this lecture)

# B+ Tree Basics

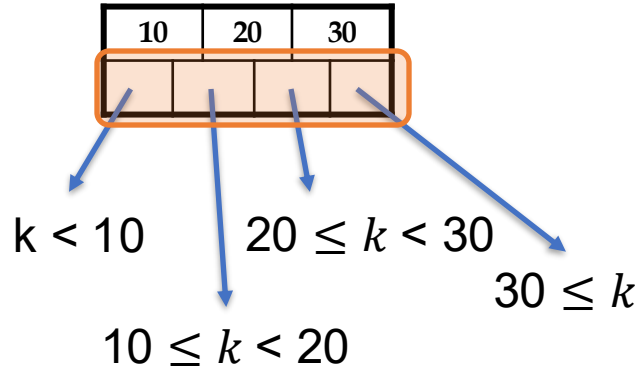


Parameter  $d$  = the degree

Each non-leaf (“interior”) node has  $\geq d$  and  $\leq 2d$  keys\*

\*except for root node, which can have between 1 and  $2d$  keys

# B+ Tree Basics

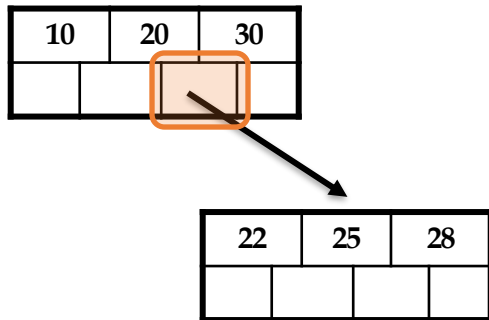


The  $n$  keys in a node define  $n+1$  ranges



# B+ Tree Basics

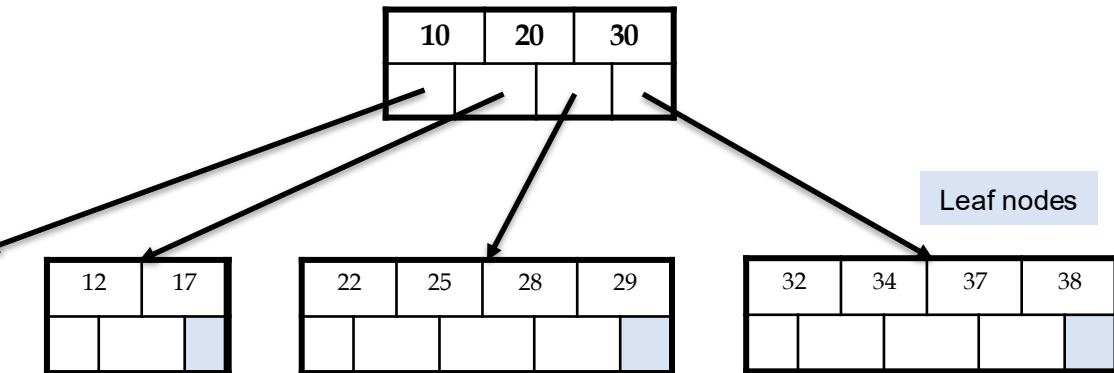
Non-leaf or *internal* node



For each range, in a non-leaf node, there is a pointer to another node with keys in that range

# B+ Tree Basics

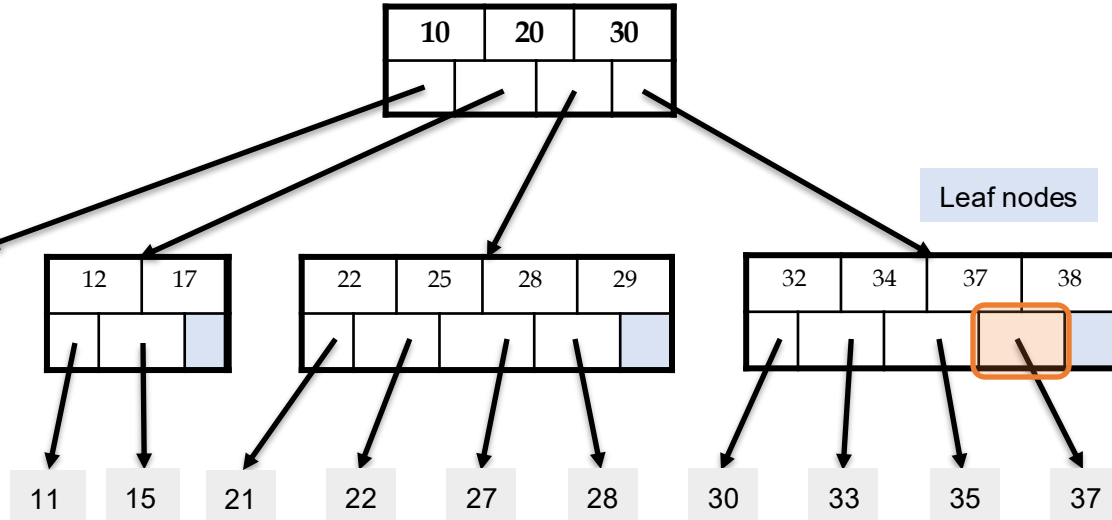
Non-leaf or *internal* node



Leaf nodes also have  
between  $d$  and  $2d$  keys,  
and are different in that:

# B+ Tree Basics

Non-leaf or *internal* node

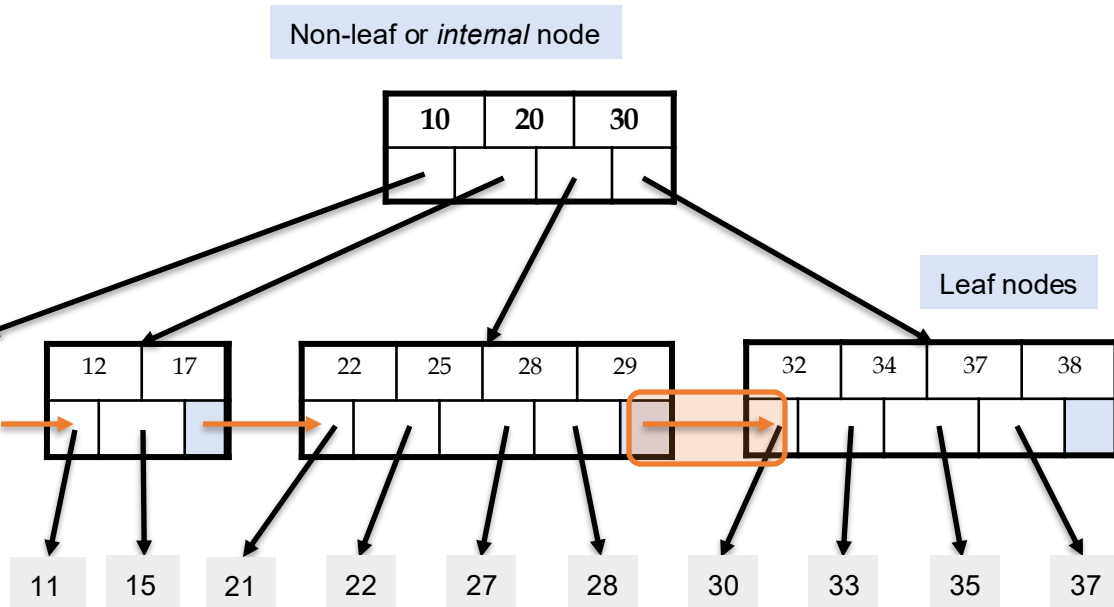


Leaf nodes

Leaf nodes also have between  $d$  and  $2d$  keys, and are different in that:

Their key slots contain pointers to data records

# B+ Tree Basics



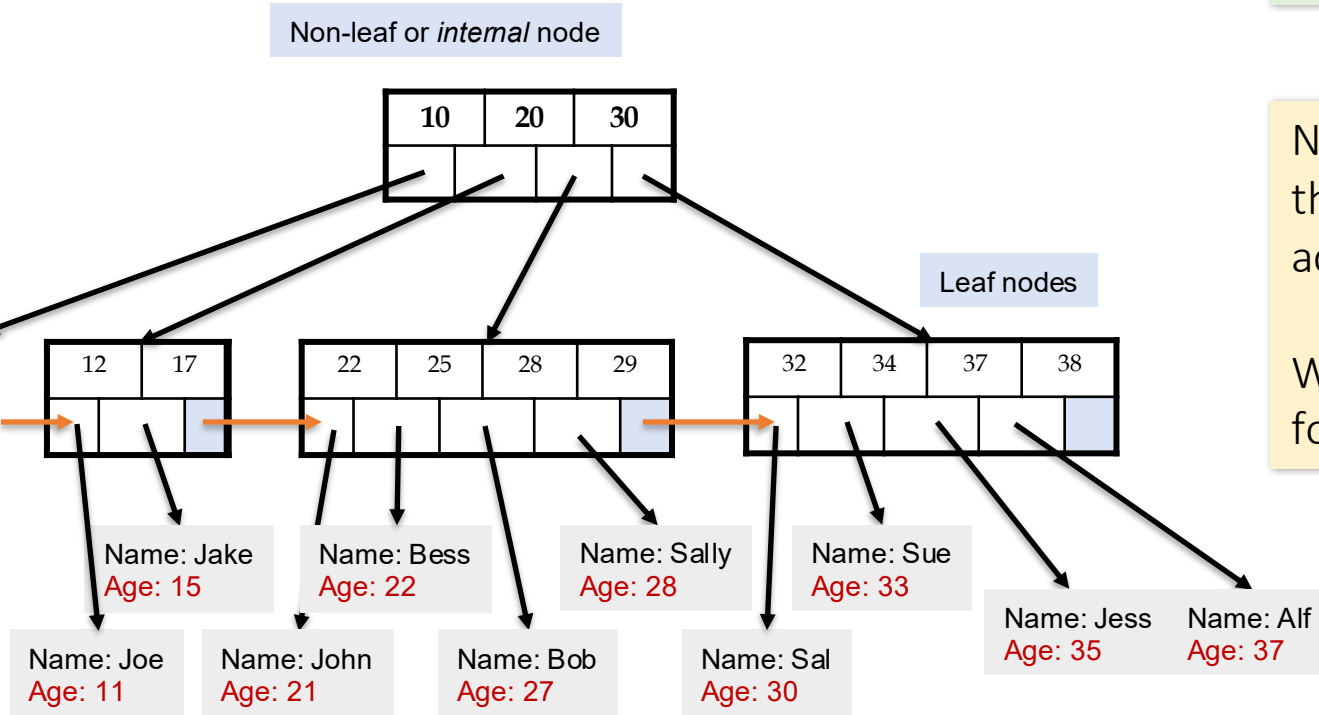
Leaf nodes also have between  $d$  and  $2d$  keys, and are different in that:

Their key slots contain pointers to data records

They contain a pointer to the next leaf node as well, for faster sequential traversal

# B+ Tree Basics

Assuming unclustered index



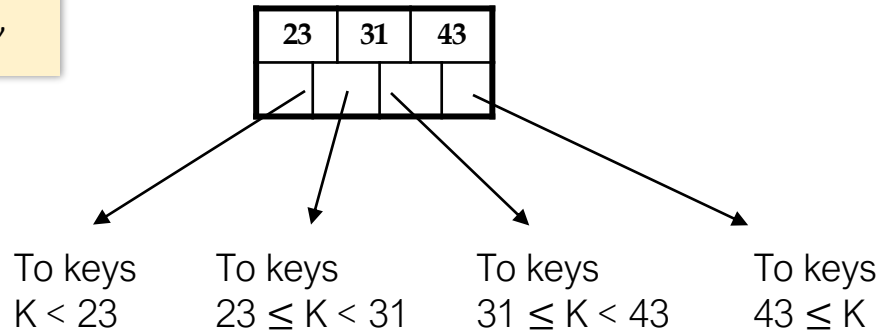
Note that the pointers at the leaf level will be to the actual data records (rows).

We might truncate these for simpler display...

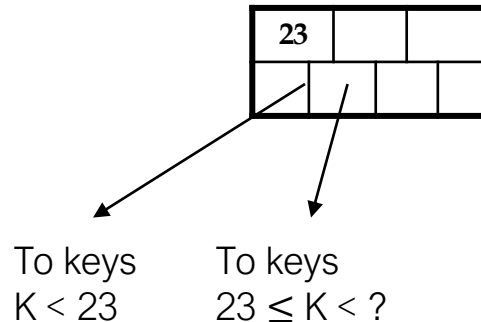
# B+ Tree occupancy requirement: interior nodes

Every node (except root) must be at least "half-full"

Full



Minimal

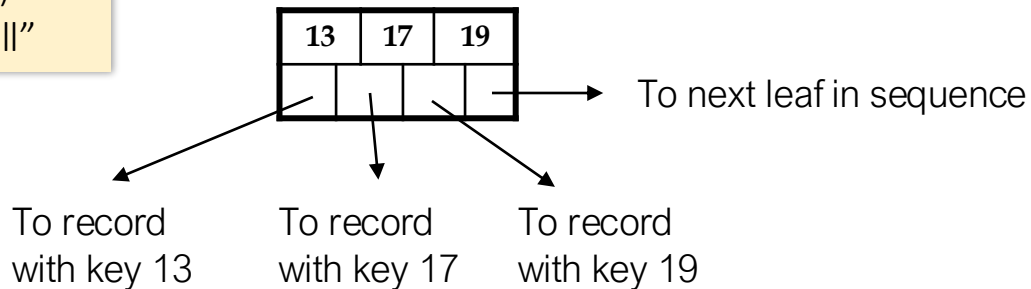


At least half of the **pointers** must be used

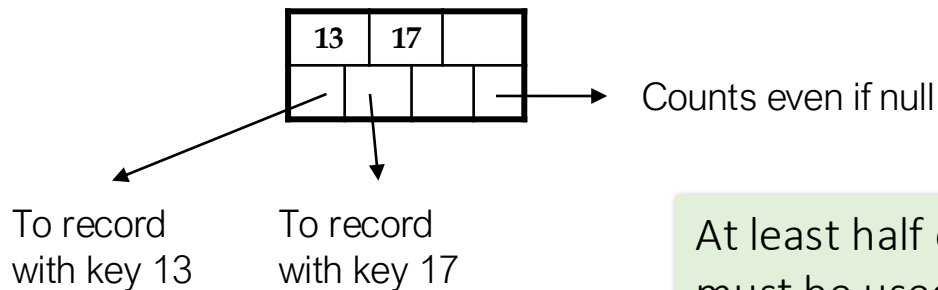
# B+ Tree occupancy requirement: leaf nodes

Every node (except root)  
must be at least "half-full"

Full



Minimal



At least half of the **keys**  
must be used

# Occupancy requirement: why does it matter?

## Ensure that the tree remains balanced

- Nodes split when they get too full
- Nodes merge/redistribute when they get too empty

Expected tree height:  $O(\log N)$   
Important for minimizing disk I/Os

## Stable and predictable performance

- Under modifications
- Also good for query planning

## Efficient space utilization

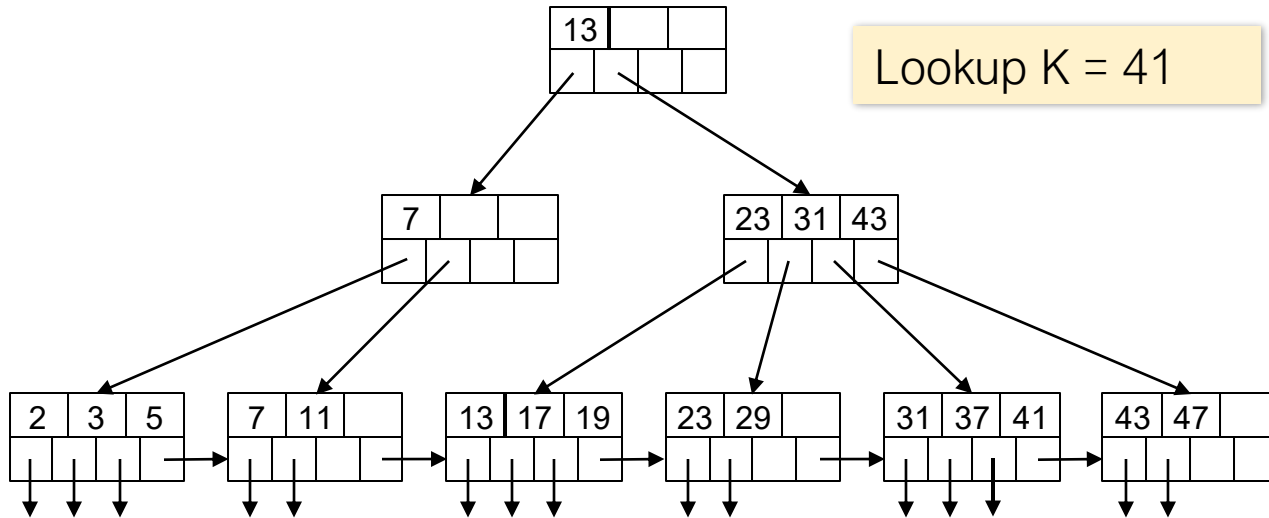
- at least 50% full



## 2. B+-Tree Operations

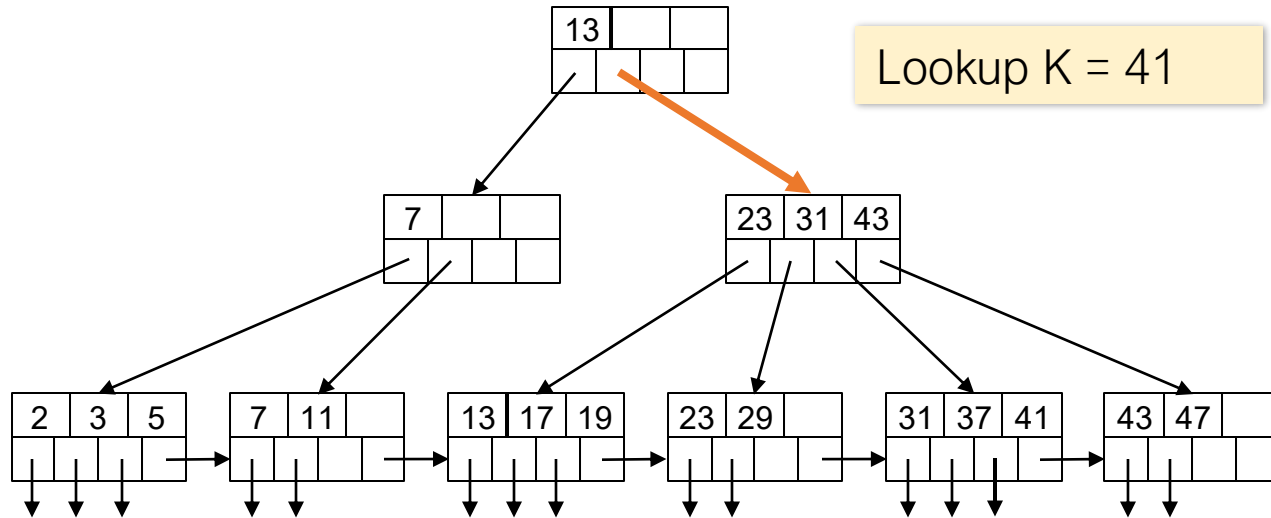
# B+ Tree: Lookup

- Search for key K recursively



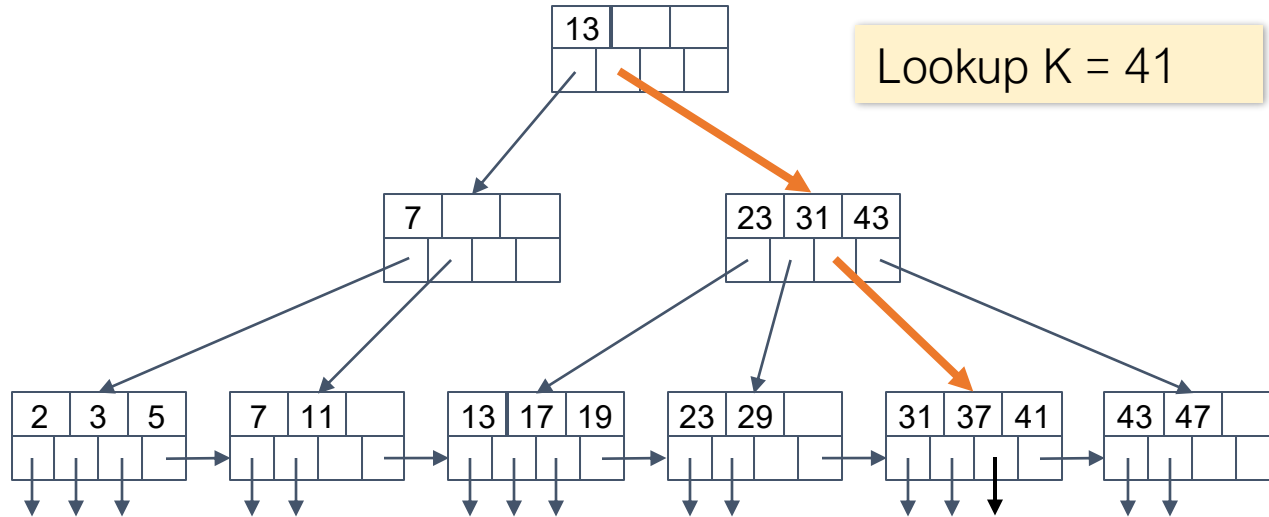
# B+ Tree: Lookup

- Search for key K recursively



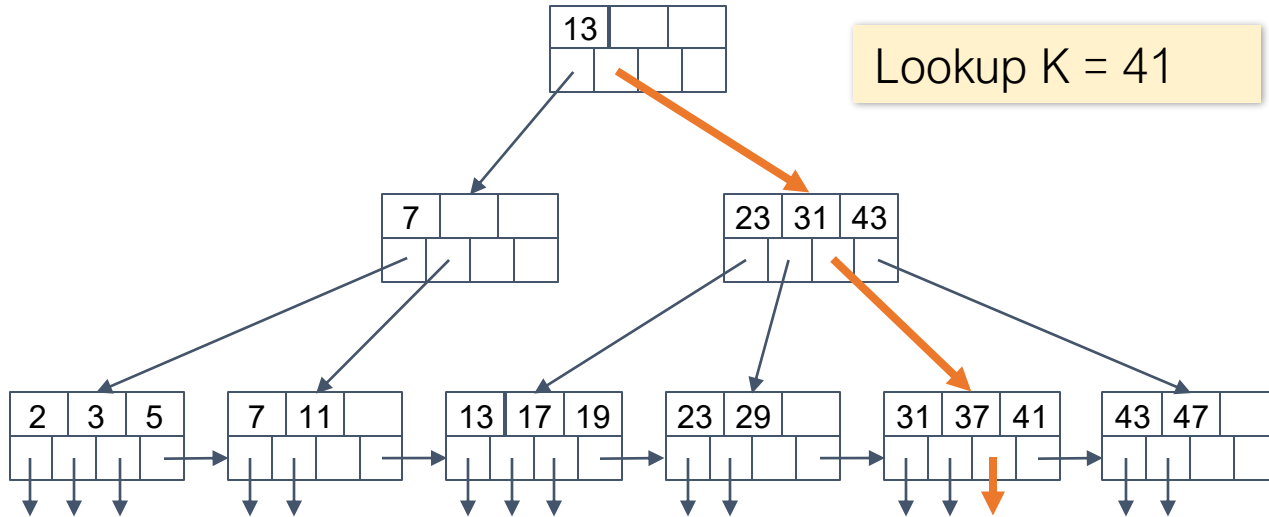
# B+ Tree: Lookup

- Search for key K recursively



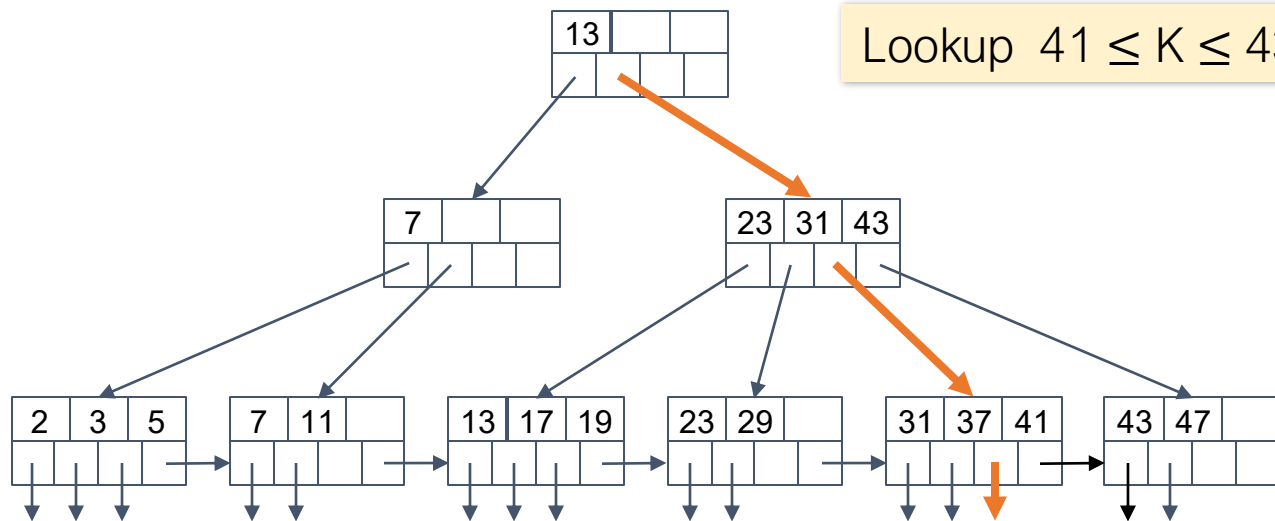
# B+ Tree: Lookup

- Search for key K recursively



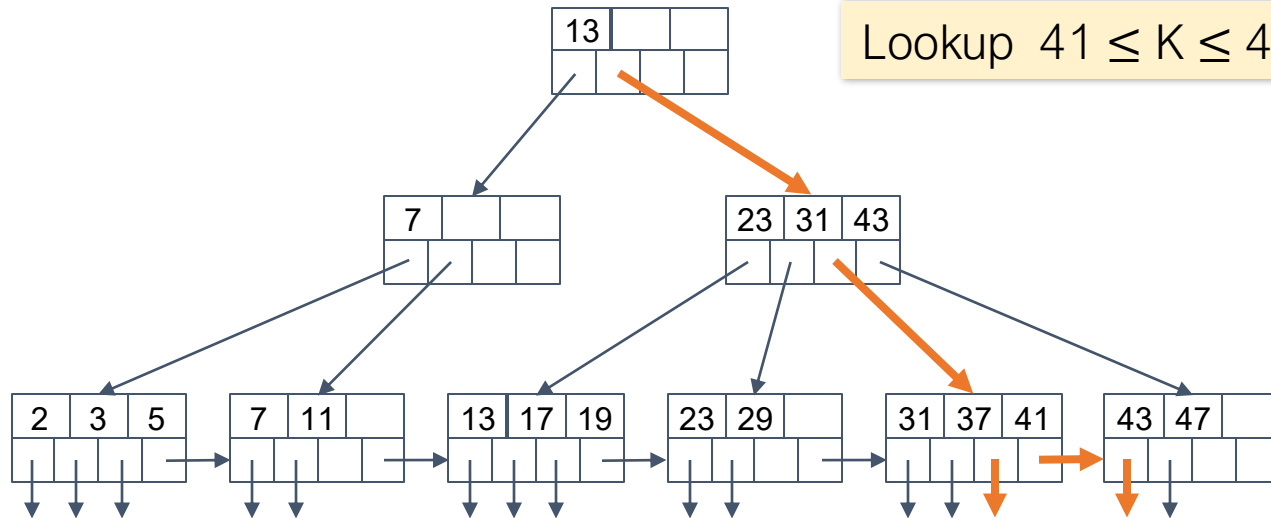
# B+ Tree: Lookup

- For range query  $[a, b]$ , search for key  $a$
- Then scan leaves to right until we pass  $b$



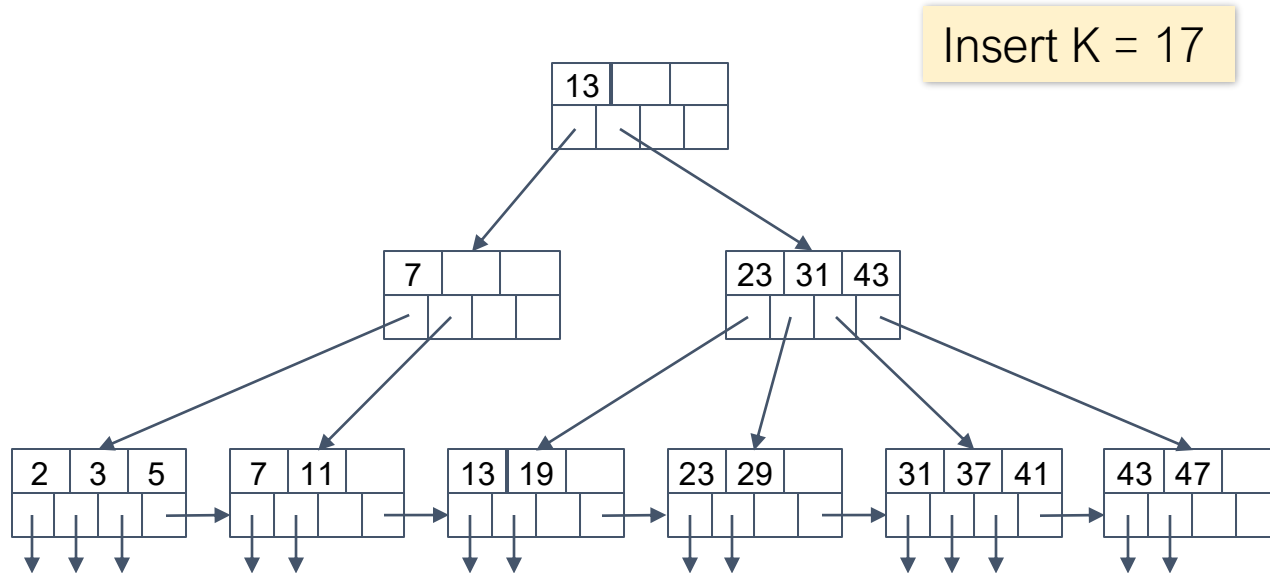
# B+ Tree: Lookup

- For range query  $[a, b]$ , search for key  $a$
- Then scan leaves to right until we pass  $b$



# B+ Tree: Insertion

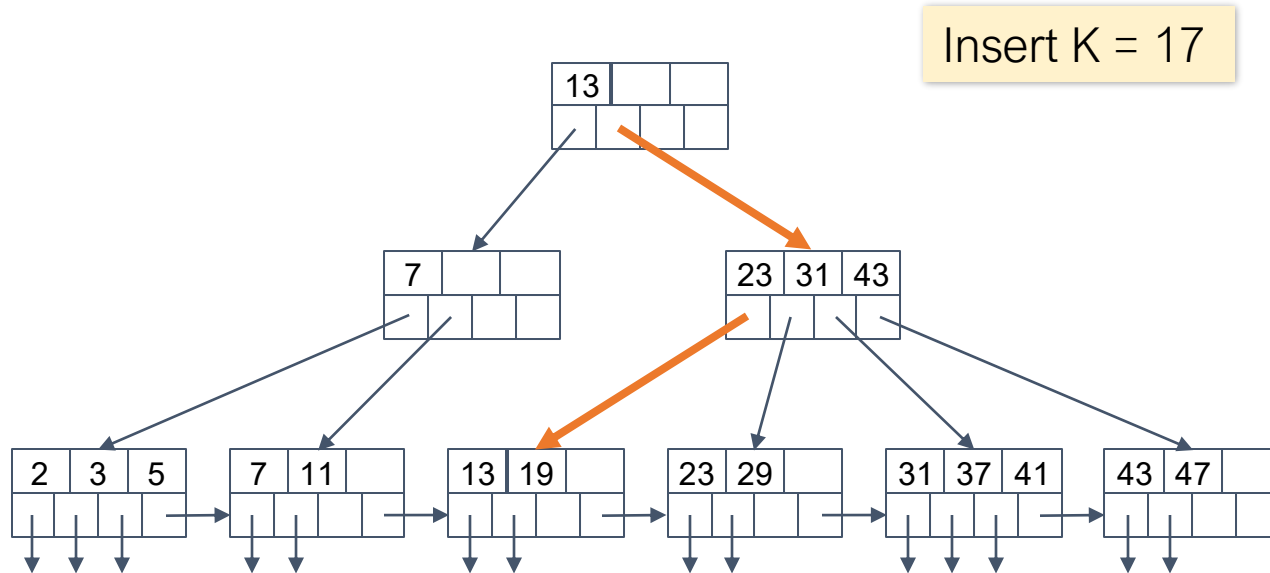
- Find place for new key in a leaf
- If there is space, put key in leaf





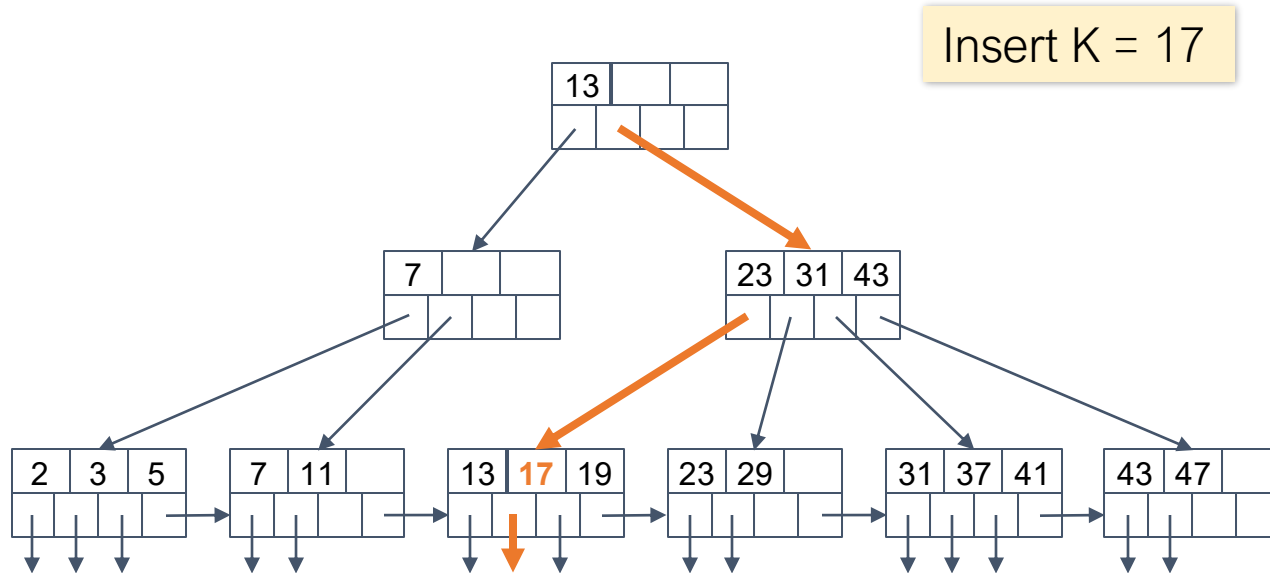
# B+ Tree: Insertion

- Find place for new key in a leaf
- If there is space, put key in leaf



# B+ Tree: Insertion

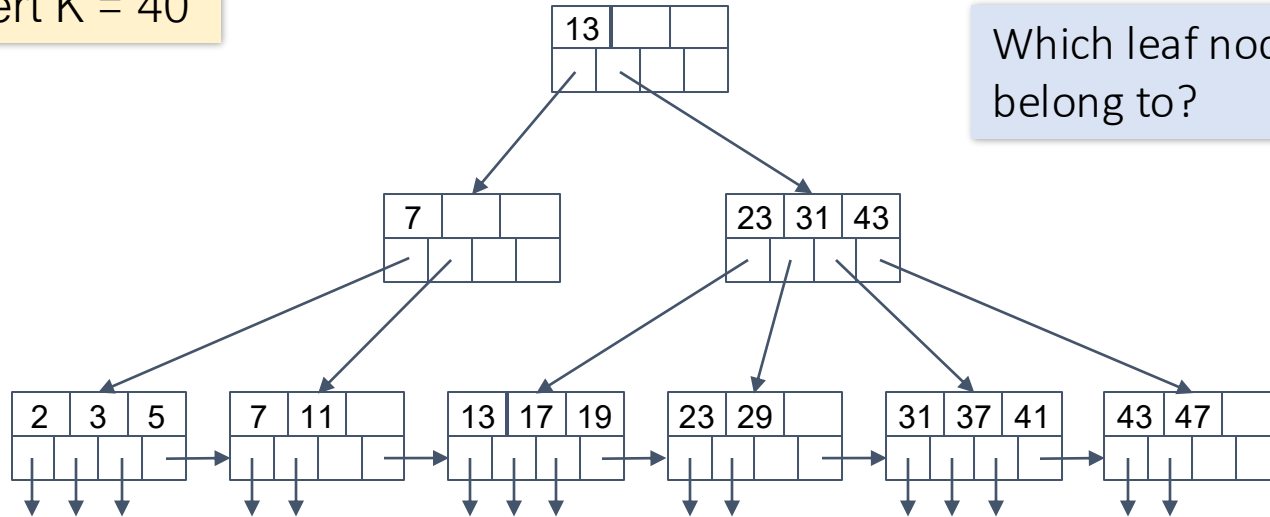
- Find place for new key in a leaf
- If there is space, put key in leaf



# B+ Tree: Insertion

- A more complex insertion example:

Insert K = 40

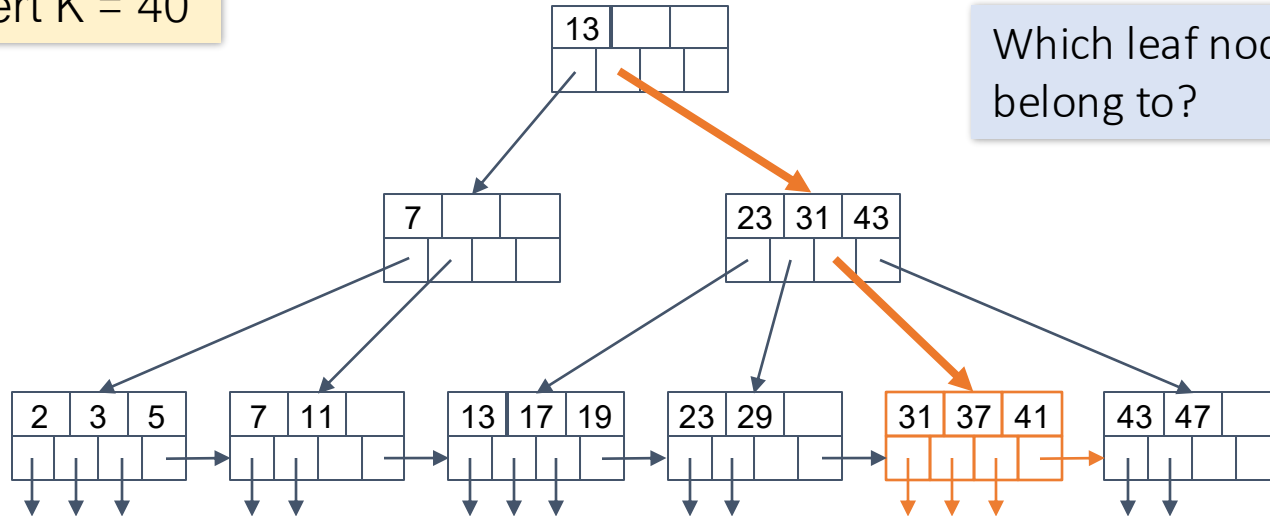


Which leaf node does K belong to?

# B+ Tree: Insertion

- If leaf is full, split into two and insert new pointer at a higher level recursively

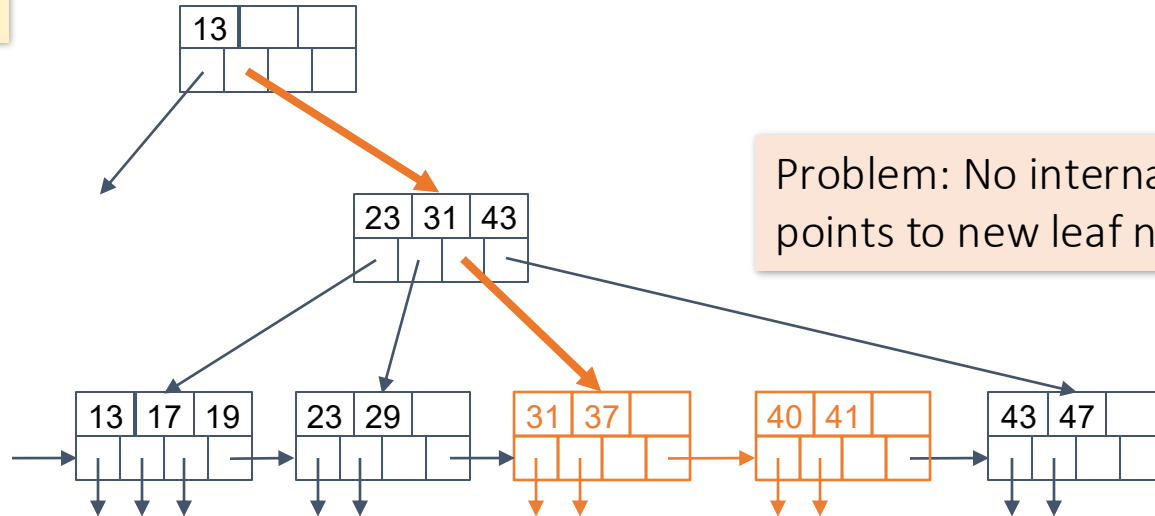
Insert K = 40



# B+ Tree: Insertion

- If leaf is full, split into two and insert new pointer at a higher level recursively

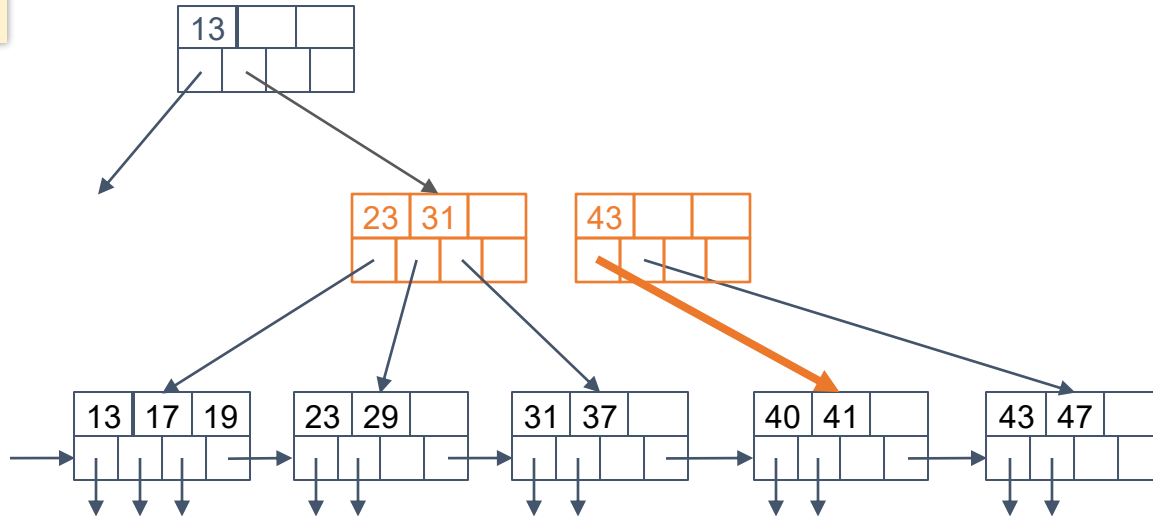
Insert K = 40



# B+ Tree: Insertion

- If leaf is full, split into two and insert new pointer at a higher level recursively

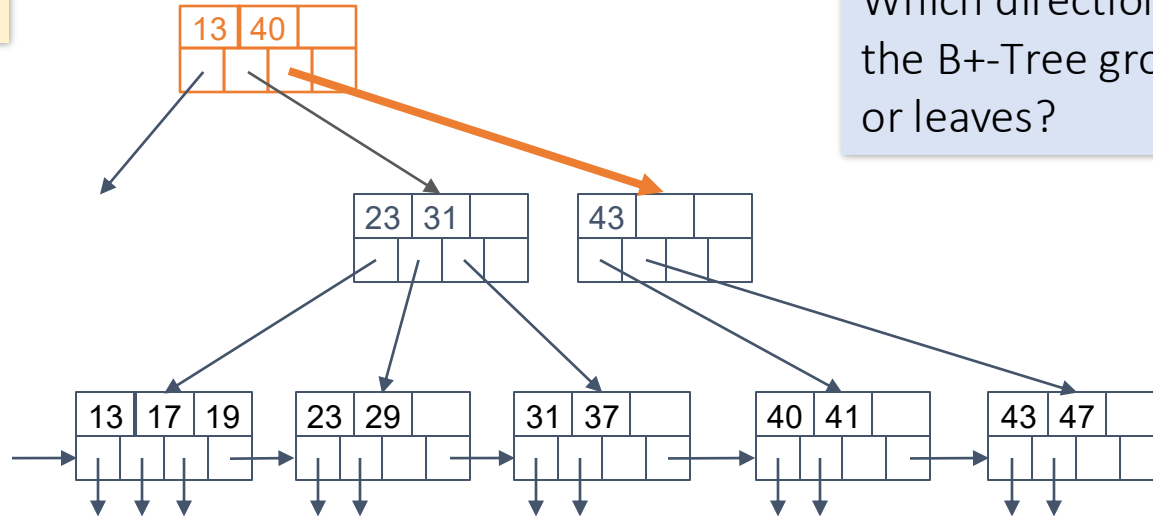
Insert K = 40



# B+ Tree: Insertion

- If leaf is full, split into two and insert new pointer at a higher level recursively

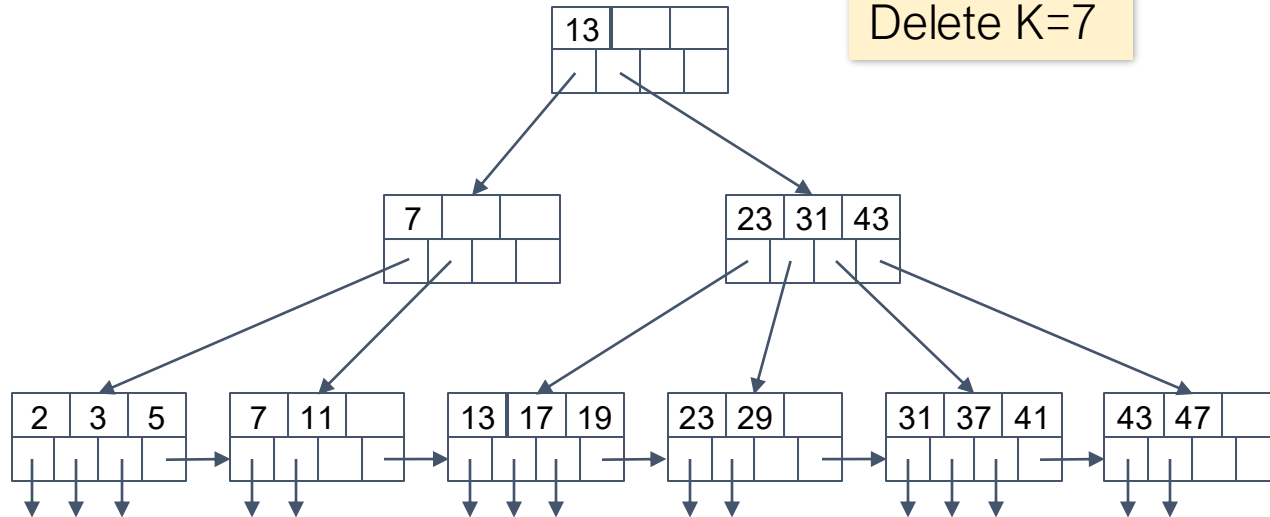
Insert K = 40



Which direction does the B+-Tree grow? Root or leaves?

# B+ Tree: Deletion

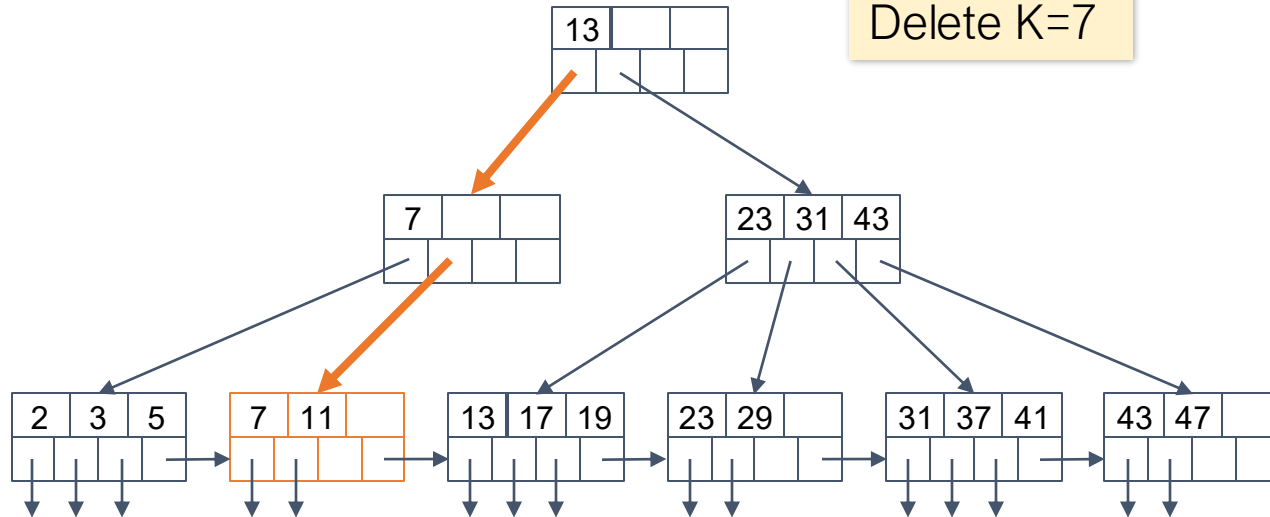
- Delete the key pointer from a leaf
- If the node contains too few pointers, take a pointer from or merge with adjacent sibling





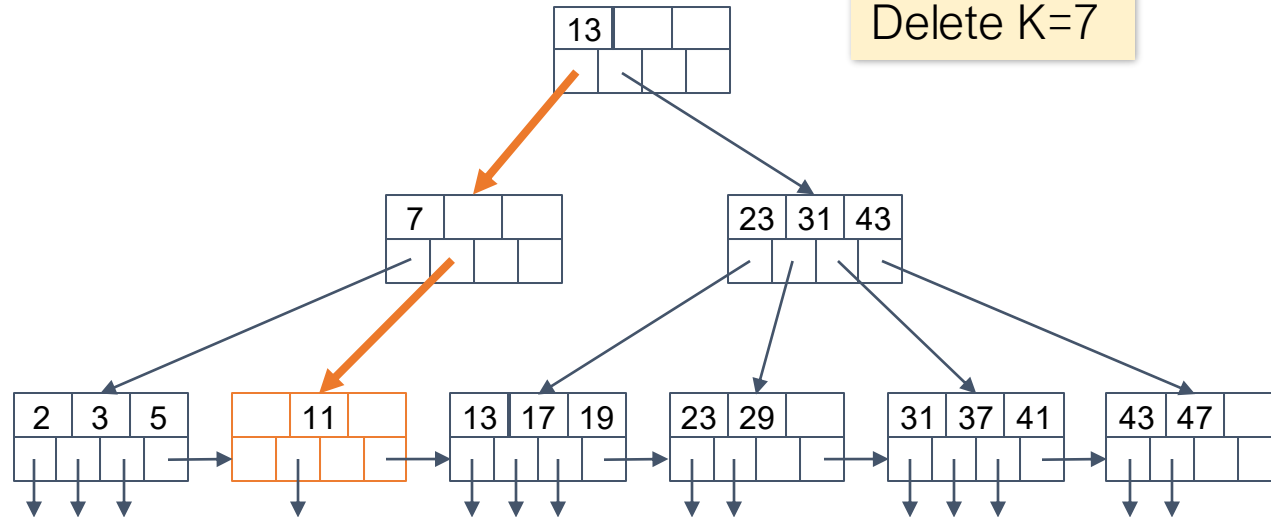
# B+ Tree: Deletion

- Delete the key pointer from a leaf
- If the node contains too few pointers, take a pointer from or merge with adjacent sibling



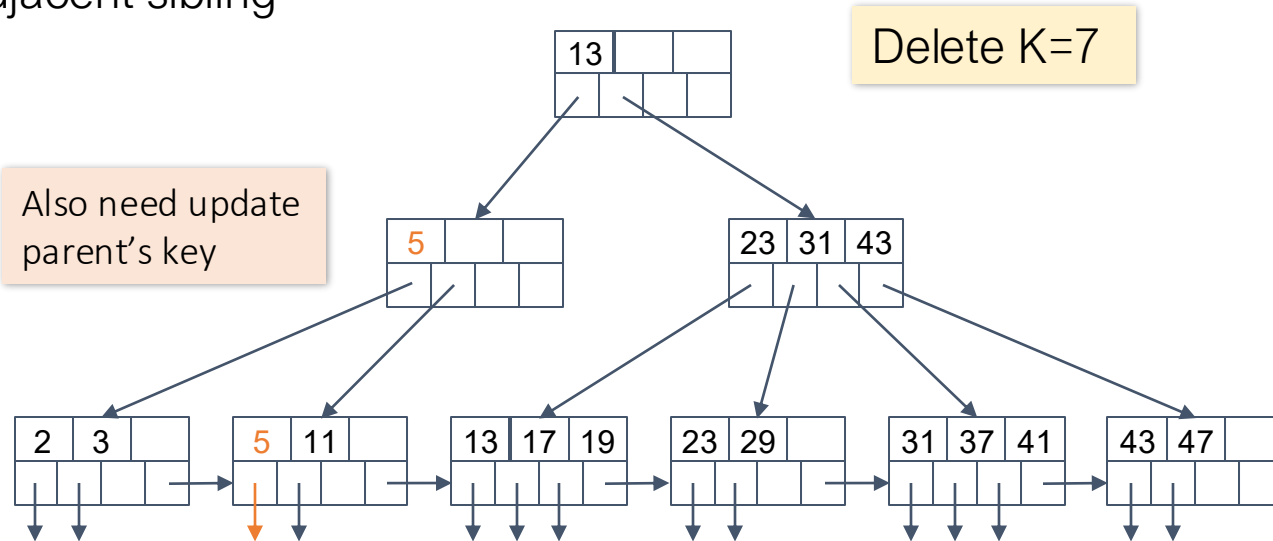
# B+ Tree: Deletion

- Delete the key pointer from a leaf
- If the node contains too few pointers, take a pointer from or merge with adjacent sibling



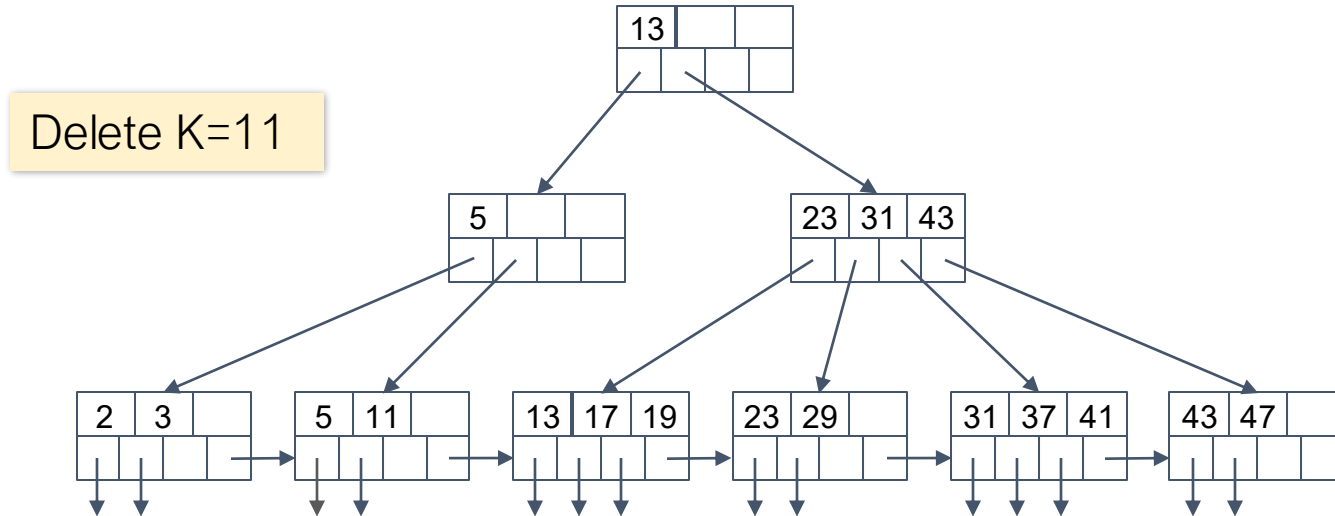
# B+ Tree: Deletion

- Delete the key pointer from a leaf
- If the node contains too few pointers, take a pointer from or merge with adjacent sibling



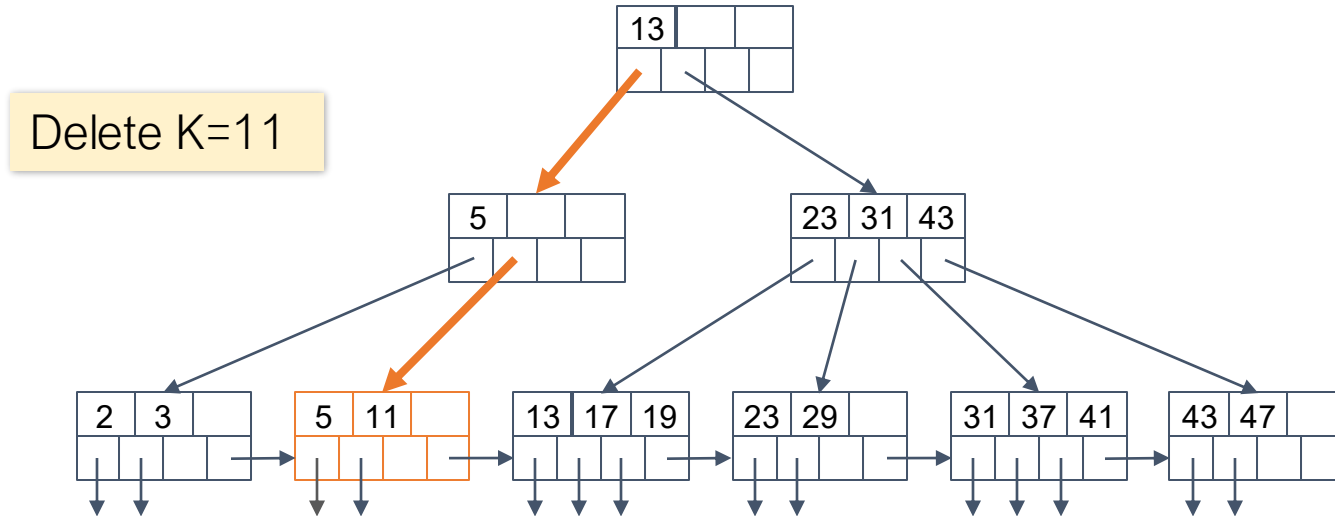
# B+ Tree: Deletion

- A more complex deletion example:



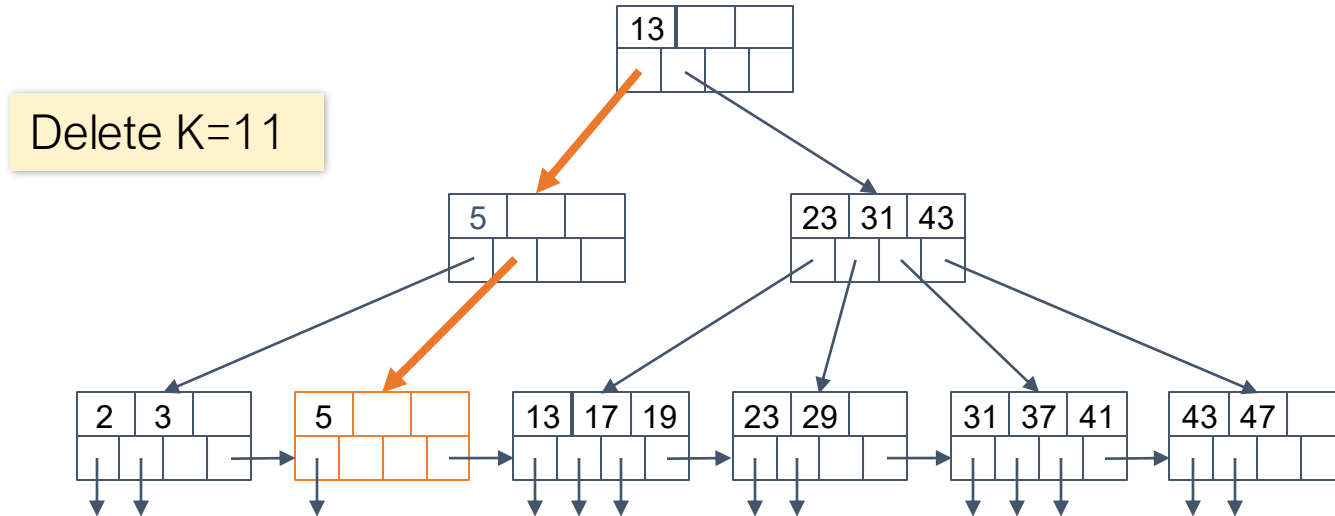
# B+ Tree: Deletion

- A more complex deletion example:



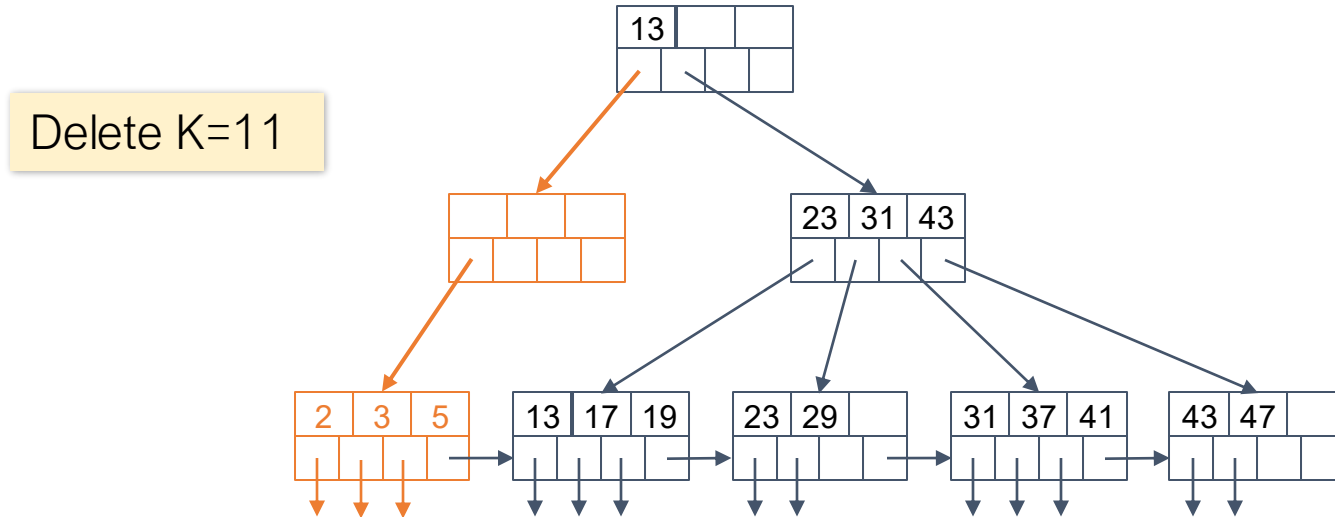
# B+ Tree: Deletion

- A more complex deletion example:



# B+ Tree: Deletion

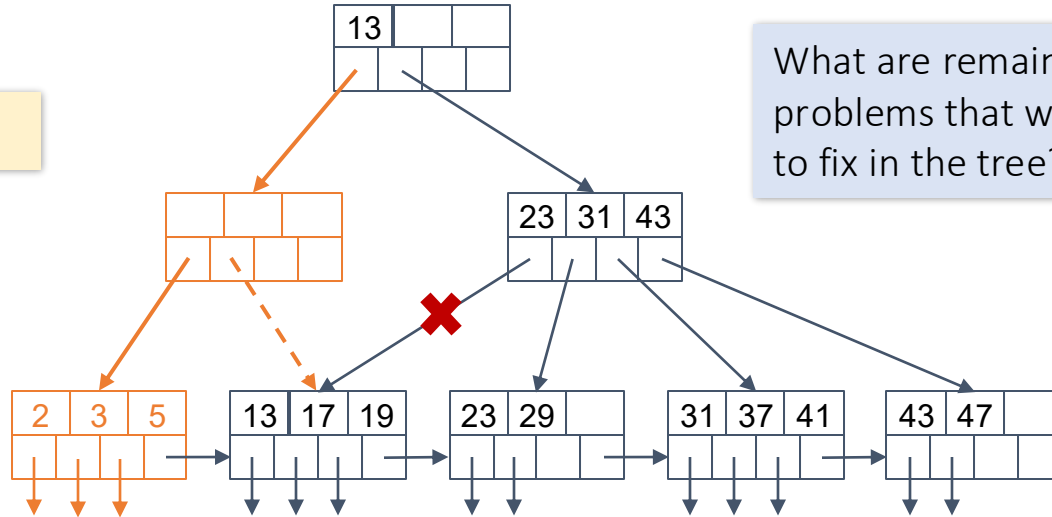
- A more complex deletion example:



# B+ Tree: Deletion

- A more complex deletion example:

Delete K=11

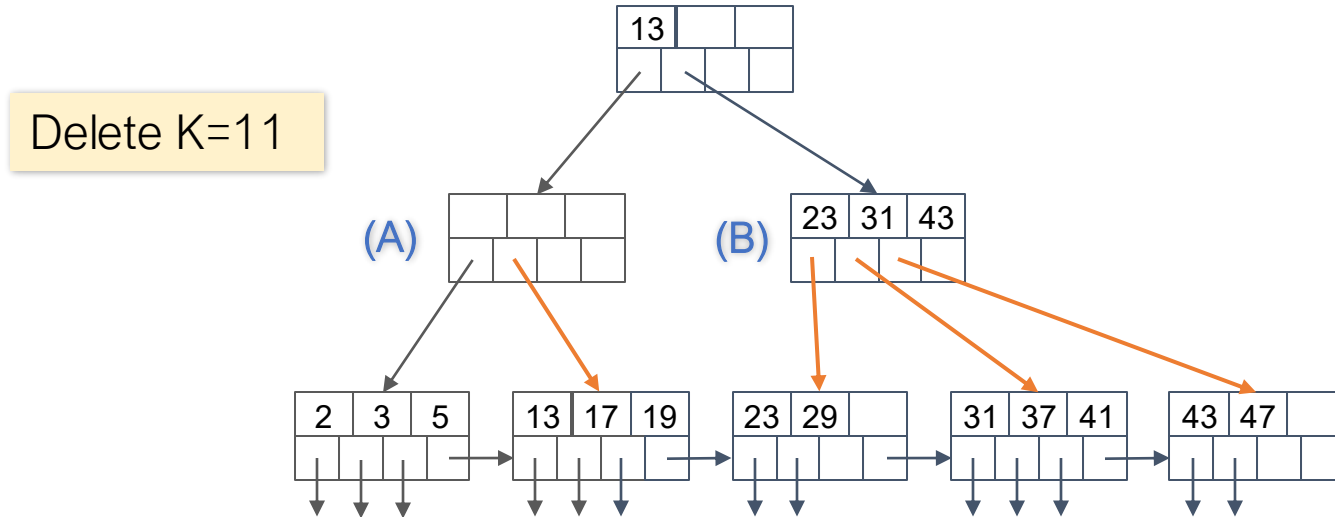


What are remaining problems that we need to fix in the tree?



# B+ Tree: Deletion

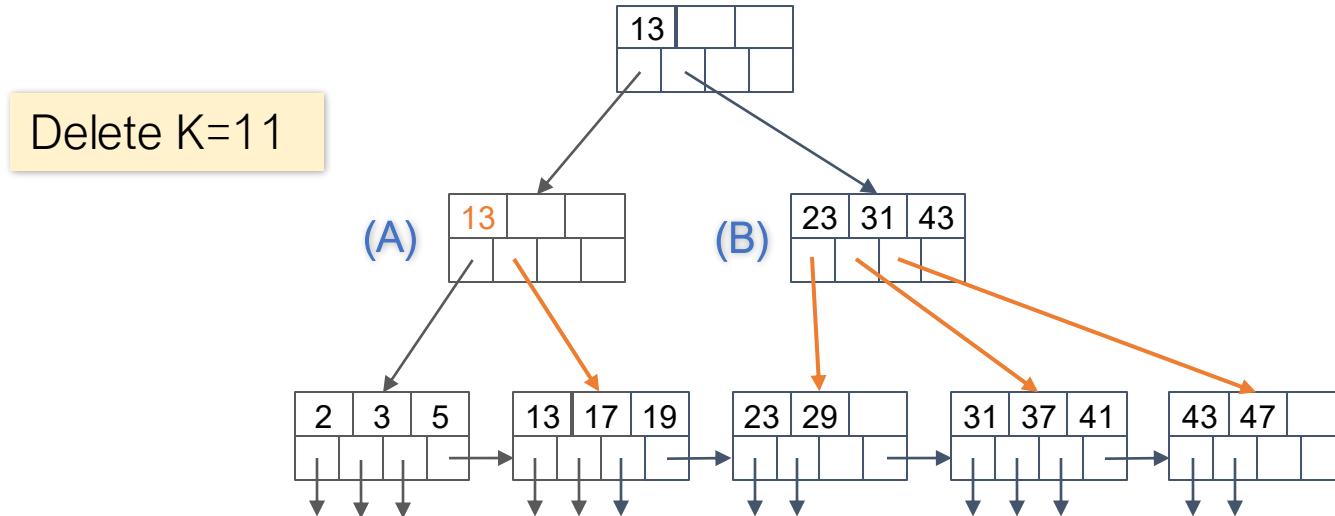
How to update keys:



# B+ Tree: Deletion

How to update keys:

- Parent key moves down to underflow node (A)



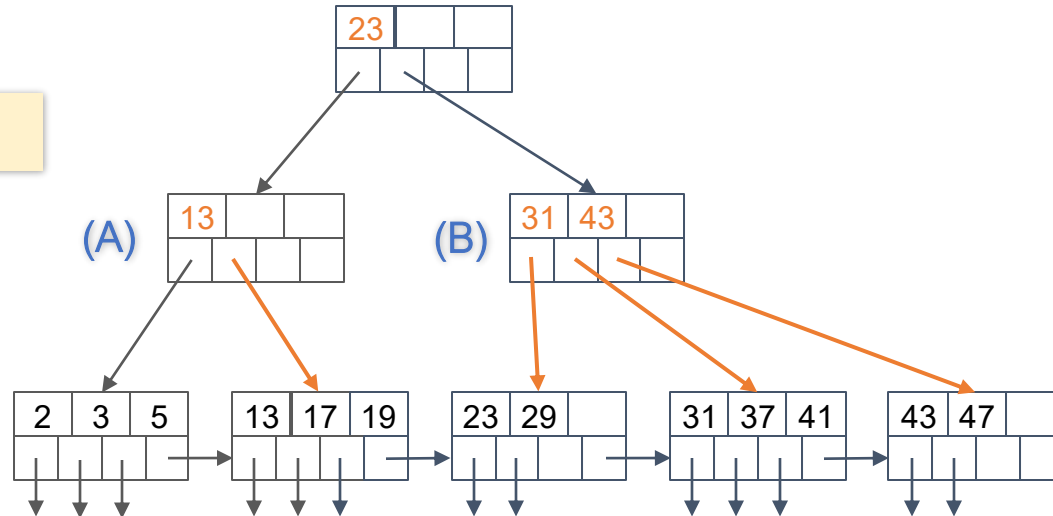
# B+ Tree: Deletion

How to update keys:

- Parent key moves down to underflow node (A)
- Smallest key from right sibling (B) moves up to parent

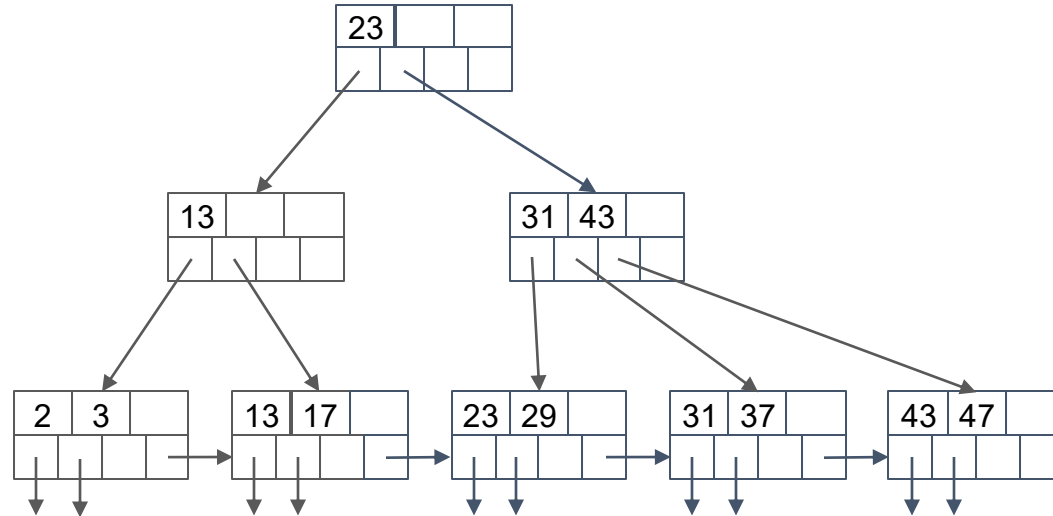
In practice, coalescing is sometimes not implemented because 1) it is hard to implement and 2) the tree will probably grow again.

Delete K=11



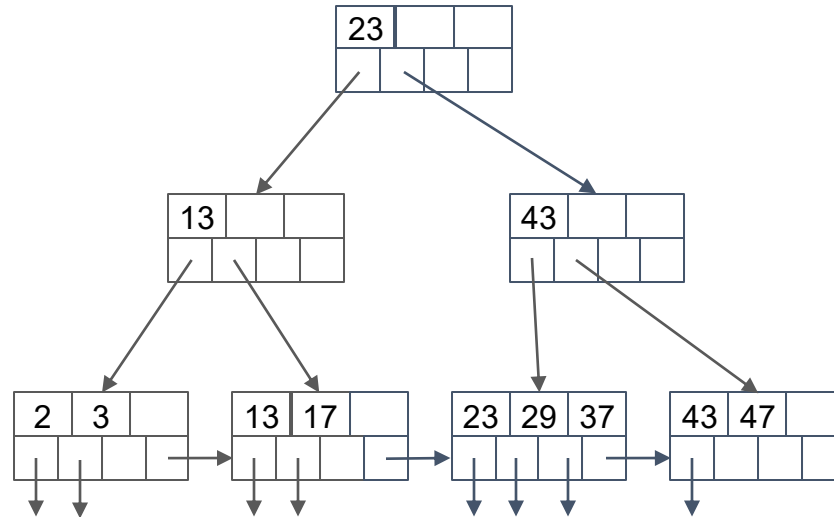
# In-class Exercise

- Delete  $K = 31$



# In-class Exercise

- Delete  $K = 31$



### 3. B+-Tree cost model

# B+ Tree: High Fanout = Smaller & Lower IO

So why does B+ tree work?

As compared to binary search trees, B+ Trees have **high *fanout*** (*between  $d+1$  and  $2d+1$* )

This means that the **depth of the tree is small** → getting to any element requires very few IO operations!

- Also can often store most or all of the B+ Tree in main memory!

The fanout is defined as the number of pointers to child nodes coming out of a node

Note that fanout is dynamic - we'll often assume it's constant just to come up with approximate eqns!

# B+ Trees in Practice

Typical order:  $d=100$ . Typical fill-factor: 67%.

- average fanout = 133

Top levels of tree sit *in the buffer pool*:

- Level 1 = 1 page = 8 KB
- Level 2 = 133 pages = 1 MB
- Level 3 = 17,689 pages = 133 MB

Fill-factor is the percent of available slots in the B+ Tree that are filled; is usually  $< 1$  to leave slack for (quicker) insertions

Typically, only pay for **one IO!**



# Simple Cost Model for Search

Suppose:

- $f$  = fanout, which is in  $[d+1, 2d+1]$  (*we'll assume it's constant for our cost model...*)
- $N$  = the total number of *pages* we need to index
- $F$  = fill-factor (usually  $\approx 2/3$ )

Our B+ Tree needs to have room to index  $N/F$  pages!

- We have the fill factor in order to leave some open slots for faster insertions

What height ( $h$ ) does our B+ Tree need to be?

- $h=1 \rightarrow$  Just the root node- room to index  $f$  pages
- $h=2 \rightarrow$   $f$  leaf nodes- room to index  $f^2$  pages
- $h=3 \rightarrow$   $f^2$  leaf nodes- room to index  $f^3$  pages
- ...
- $h \rightarrow f^{h-1}$  leaf nodes- room to index  $f^h$  pages!

$\rightarrow$  We need a B+ Tree of height  $h = \left\lceil \log_f \frac{N}{F} \right\rceil$ !

# Simple Cost Model for Search

Note that if we have  $B$  available buffer pages, by the same logic:

- We can store  $L_B$  levels of the B+ Tree in memory
- where  $L_B$  is the number of levels such that the sum of all the levels' nodes fit in the buffer:
  - $B \geq 1 + f + \dots + f^{L_B-1} = \sum_{l=0}^{L_B-1} f^l$

In summary: to do exact search:

- We read in one page per level of the tree
- However, levels that we can fit in buffer are free!
- Finally we read in the actual record

$$\text{IO Cost: } \left\lceil \log_f \frac{N}{F} \right\rceil - L_B + 1$$

$$\text{where } B \geq \sum_{l=0}^{L_B-1} f^l$$

# Simple Cost Model for Search

To do range search, we just follow the horizontal pointers

The IO cost is that of loading additional leaf nodes we need to access + the IO cost of loading each *page* of the results- we phrase this as “Cost(OUT)”

$$\text{IO Cost: } \left\lceil \log_f \frac{N}{F} \right\rceil - L_B + \text{Cost}(\text{OUT})$$

$$\text{where } B \geq \sum_{l=0}^{L_B-1} f^l$$

# In-class Exercise

Given a B+ tree indexing over  $N = 100,000$  data pages with fill factor  $F = 1$  and constant fanout  $f = 10$ . Assume that each node of the B+ tree occupies one page, that there are  $B = 11$  buffer pages available to store B+ Tree nodes.

What's IO cost of performing an exact search query on this index?

$$\text{IO Cost: } \left\lceil \log_f \frac{N}{F} \right\rceil - L_B + 1$$

$$\text{where } B \geq \sum_{l=0}^{L_B-1} f^l$$