

# Tree-Structured Indexes

## Chapter 10

# Recap – Storage and Indexing

- Can store a database table in
  - Heap: i.e., no indexing. (Not same as the Heap data structure in EECS 281).
  - Indexed. *Create index*
    - Searching: Given a search key  $k$ , look up matching  $K^*$  entries, Retrieve record(s) from those entries.
  - Three alternatives for  $K^*$ : *Alternative 1 ↓ = record order entry order*
    1.  $(k, \text{Record})$ . This can be considered to be a clustered index.
    2.  $(k, \text{RID})$ .  $\text{RID}$  points to a block/offset/size of  $\text{Record}$  on the disk
    3.  $(k, [n, [\text{RID1}, \text{RID2}, \dots \text{RIDn}]])$

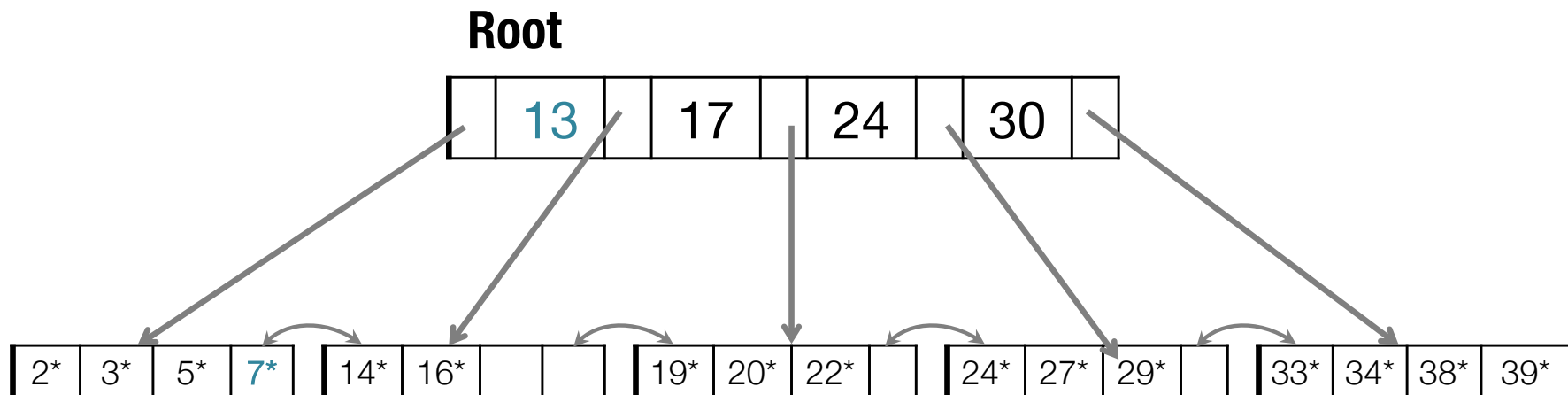
# This lecture

- Range and equality searches very common
- Can scan heap file, but this is expensive
- **Goal:** Create a dynamic index structure that allows for efficient evaluation, and equality / range queries
- B+-tree: A search tree but persistent with each tree node stored in a disk block

# Terminology

- First there was a B Tree.
- People immediately loved it
  - Many people started making improvements to it.
  - B+ Tree, B\* Tree, ...
- Eventually, all agreed that B+ Tree was the best of these data structures in the BTree family.
- Name “B+ Tree” is a twister
  - Often casually referred to as BTree.

# Example B+-tree of height 1



- Height-balanced (dynamic) search tree
- Persistent: Each node in the tree is stored in a disk block!
  - Thus, each tree node has many pointers, unlike a binary search tree!
  - Given a key, say 16, efficient to find its entry, 16\*.
  - Range queries (e.g., all records between 16-25) also can be done efficiently

# B+ Tree

- Additional Properties:

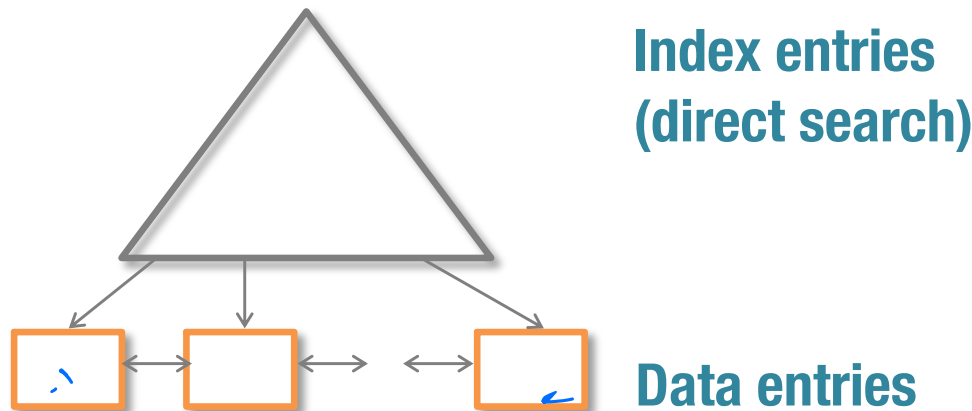
- Invariant:** Minimum 50% occupancy (except for root), i.e. each node contains  $d \leq m \leq 2d$  entries, where  $d$  = order of the tree.  
Root has min of 1 entry and max of  $2d$  entries
- Invariant:** Remains balanced after insert/delete operation

*min amount of entry in a page*

Height of tree: the length of any path from the root to any leaf.

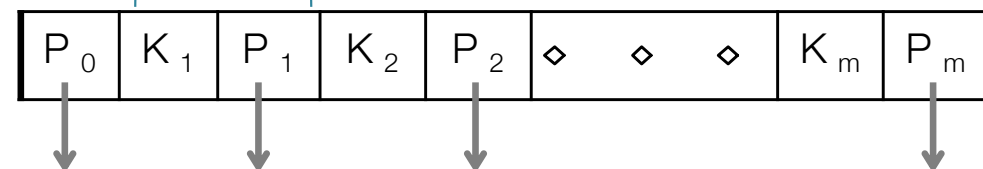
*Fan out = # children for non-leaf node*

*2d entries: dividers  
⇒ 2d + 1 children*



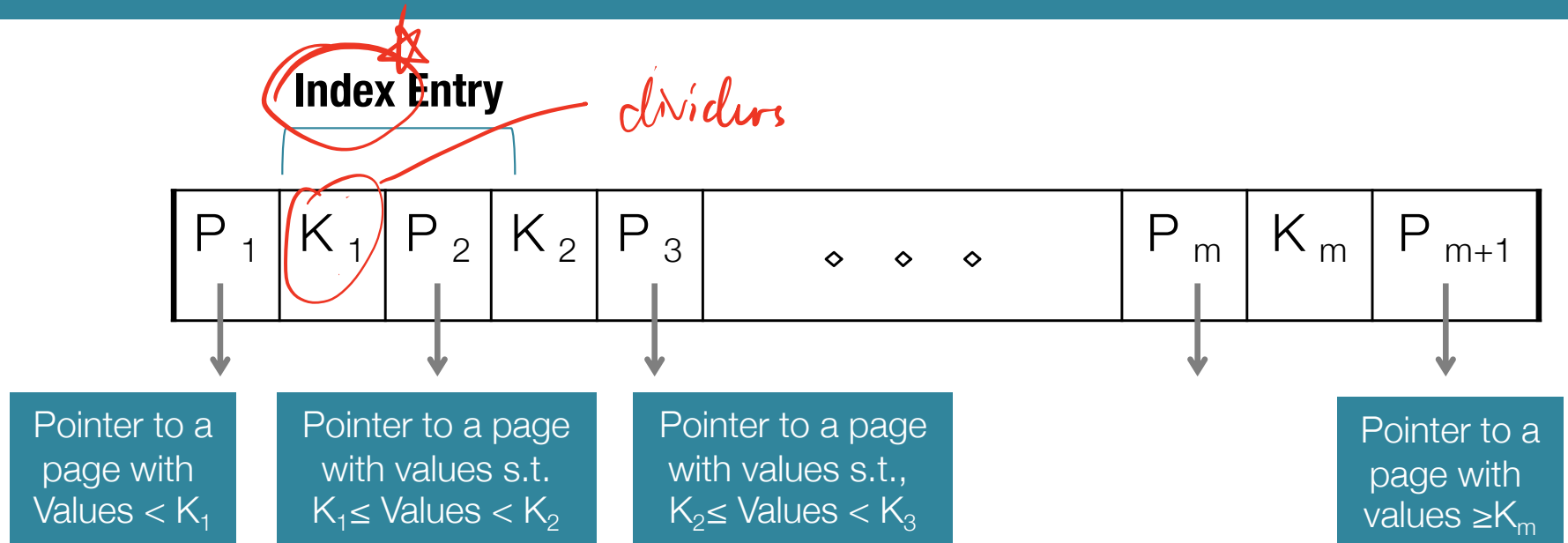
Index Entries:  
Entries in the non-leaf pages:  
(search key value, pageid)

## Index Entry

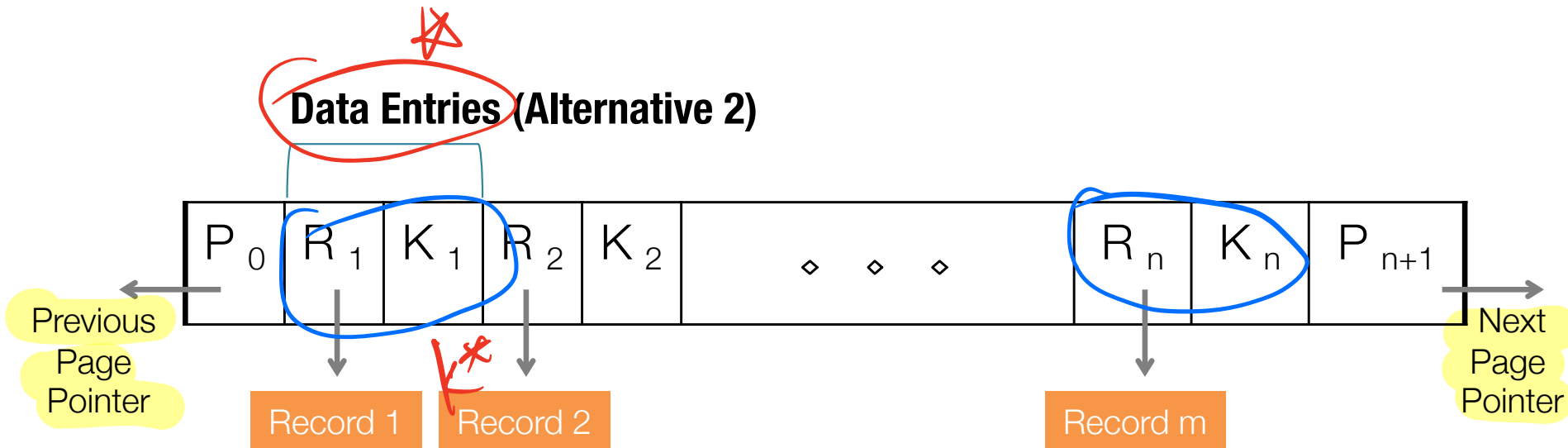


# B+-Tree Page Format

Non-leaf Page

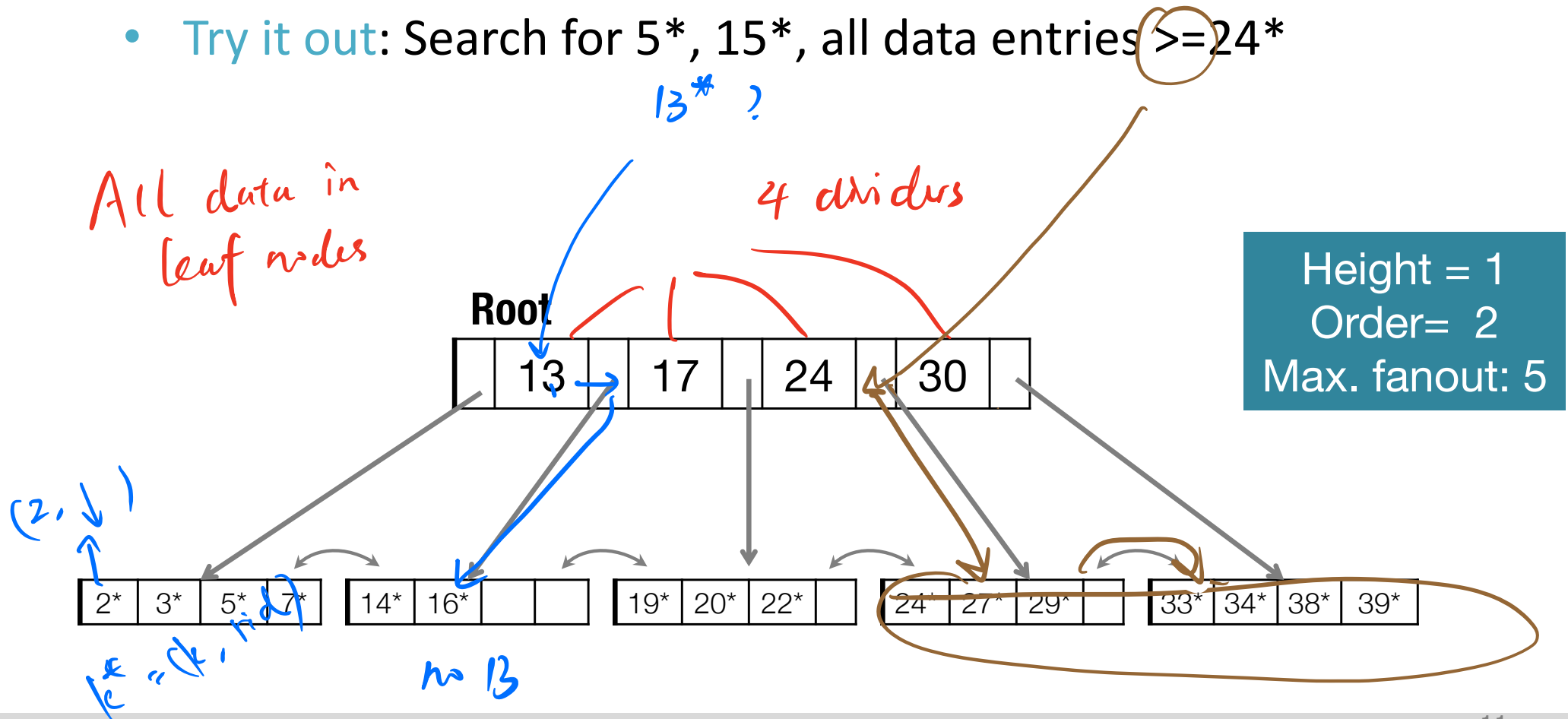


Leaf Page



# Searching in a B+ Tree

- Starting from root, examine index entries in non-leaf nodes, and traverse down the tree until a leaf node is reached
  - Non-leaf nodes can be searched using a binary or a linear search.
- Try it out: Search for 5\*, 15\*, all data entries  $\geq 24^*$





# B+ Trees: Height

What is the height of a B+ tree? (formula)

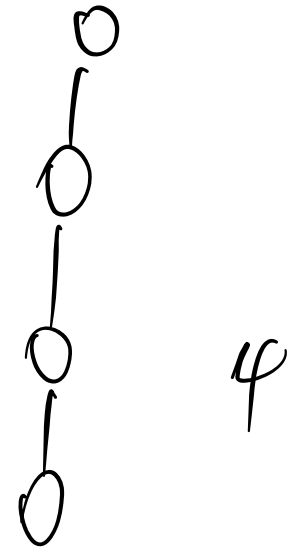
**Let:**

**F:** average fanout (average number of children for non-leaf node) *(large #)*

**N:** total number of leaf pages

*max fan out =  $2d+1$*

Then,  
height =  $\lceil \log_F(N) \rceil$



*- | - | -*

# B+ Trees in Practice

- Typical order: 100. Typical fill-factor: 2/3
  - Maximum fanout: 201  $= 2d+1$
  - Average fanout = 133
- Typical capacity:
  - Height = 1: 133 pages of data entries (leaf pages)
  - Height = 2:  $133^2$  pages of data entries
  - Height = 3:  $133^3$  (> 2 million) pages of data entries
  - Height = 4:  $133^4$  (> 300 million) pages of data entries
- Can often keep top levels of index in buffer pool
  - Level 1 = 1 page = 8 KB
  - Level 2 = 133 pages = 1 MB
  - Level 3  $\approx$  17.7K pages = 133 MB
  - Level 4  $\approx$  2.35M pages  $\approx$  17.7 GB

$$\text{height} = \log_{(2d+1) \cdot f} \left( \frac{\# \text{ records}}{(2d \cdot f)} \right)$$

# Arithmetic Example

- You are given a file of 10 million records
- Suppose you store an average of 10 data entries per leaf page
- You build a B+ Tree with order 75, 67% average fill-factor <sup>need 1 million leaf pages</sup>  
<sub>151</sub>
- What is the avg fanout? <sub>~100</sub>
- What is the height of your B+ Tree?

$$\log_{100} 1 \text{ million} = 3$$



# Question??

The average fanout is close to

C

A. 67

B. 75

C. 100

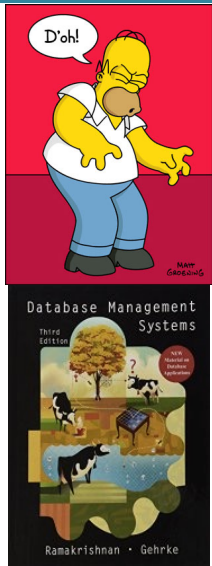
D. 1000000

# Question??

The height of the tree is *B*

- A. 2
- B. 3
- C. 4
- D. More than 4

# A Note on Order



- Some literature uses:  $\text{order} = \text{max} \# \text{ of entries}$
- In this class:  $\text{order} = \text{minimum} \# \text{ of entries}$   
(book uses this)
- *Order  $d$*  concept suggested by physical space criterion in practice (e.g. **at least half-full**)
  - Index (i.e. non-leaf) pages may hold many more entries than leaf pages, particularly for alternative 1.

# B+ Tree Operations

- Search
  - Equality
  - Range
- Insert data entry
- Delete data entry
- Bulk load

# B+ Tree: Inserting a Data Entry

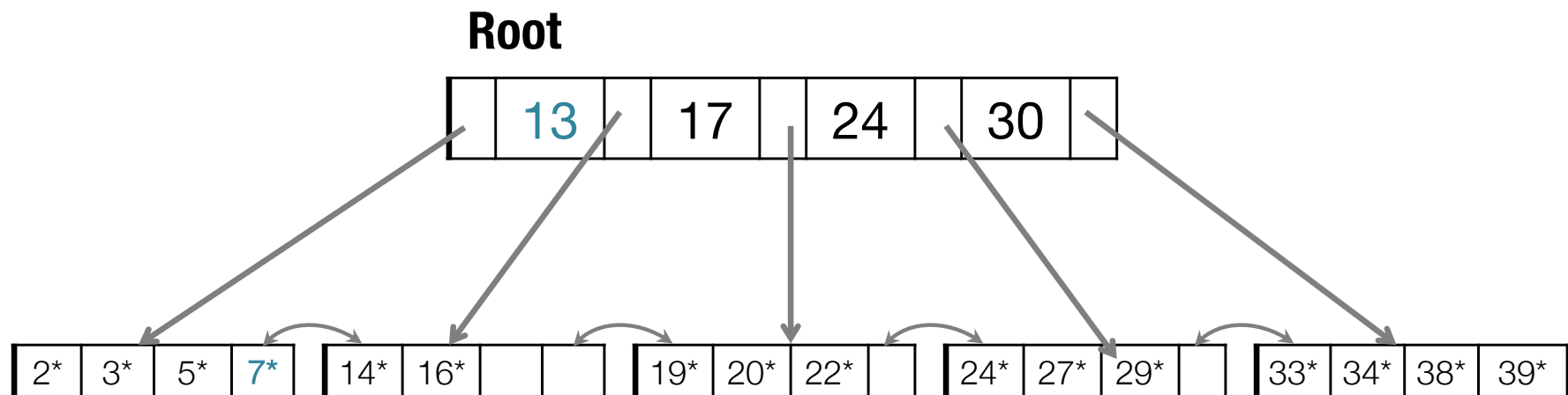


- Maintain invariants:
  - Search-tree property
  - All nodes must be at least  $\frac{1}{2}$  full (except root node), i.e., has between  $d$  and  $2d$  entries
  - Root node is allowed to have a single entry
- Strategy:
  - Split nodes when they become full and a node is added:
    - Split an overfull node with  $(2d + 1)$  entries into two nodes with  $d$  and  $(d+1)$  entries, resp.



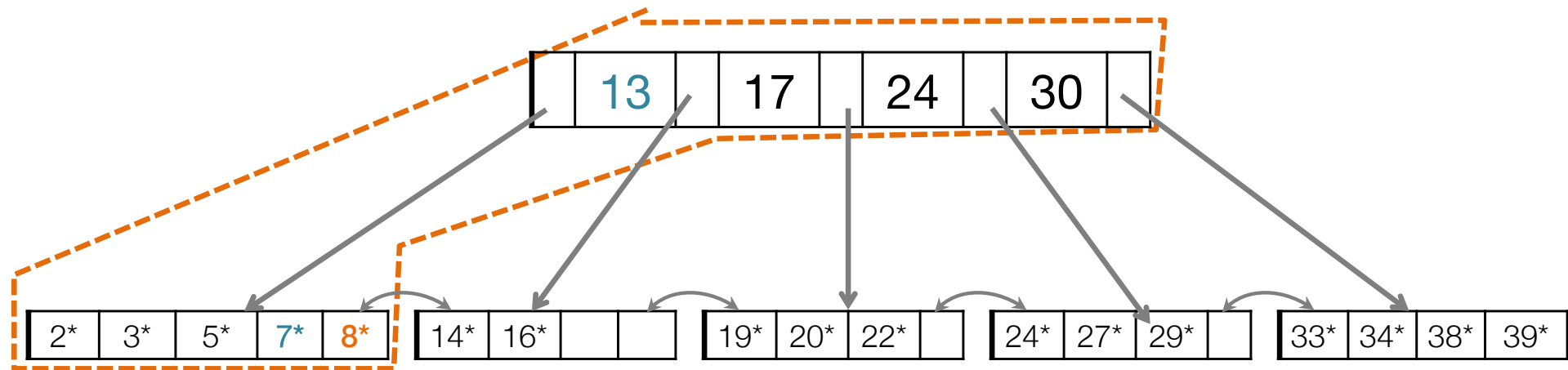


# Inserting 8\* into B+ Tree

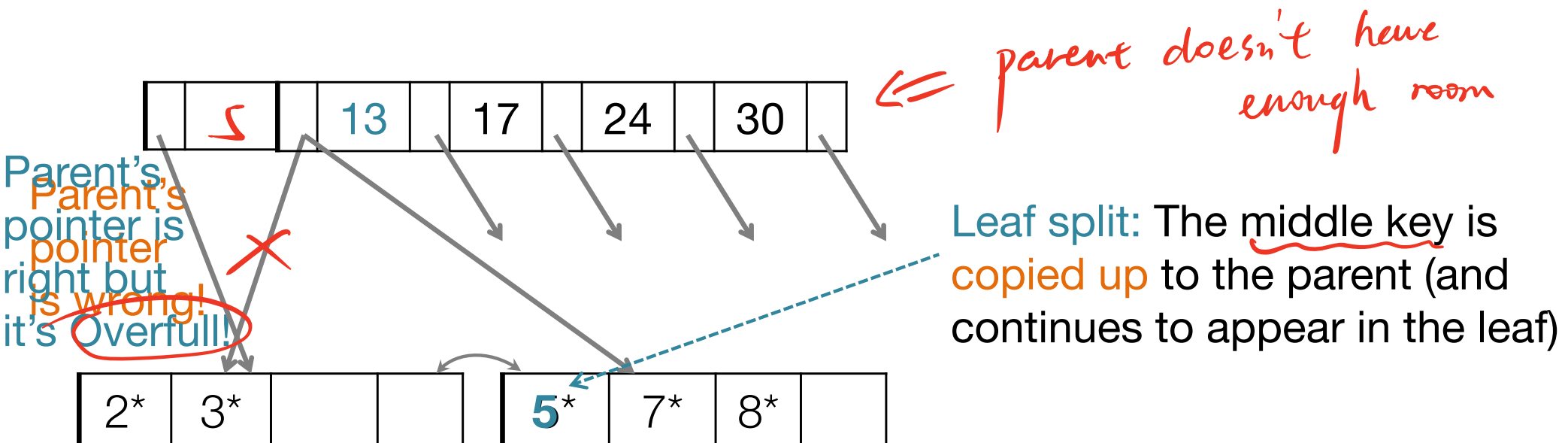


# Inserting 8\* into B+ Tree

Root

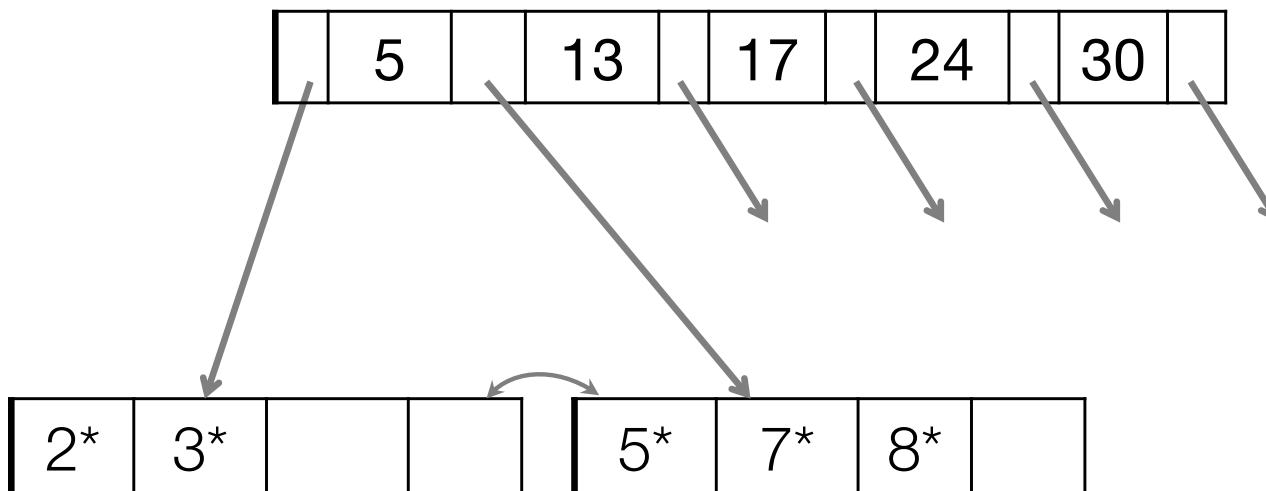


Overfull



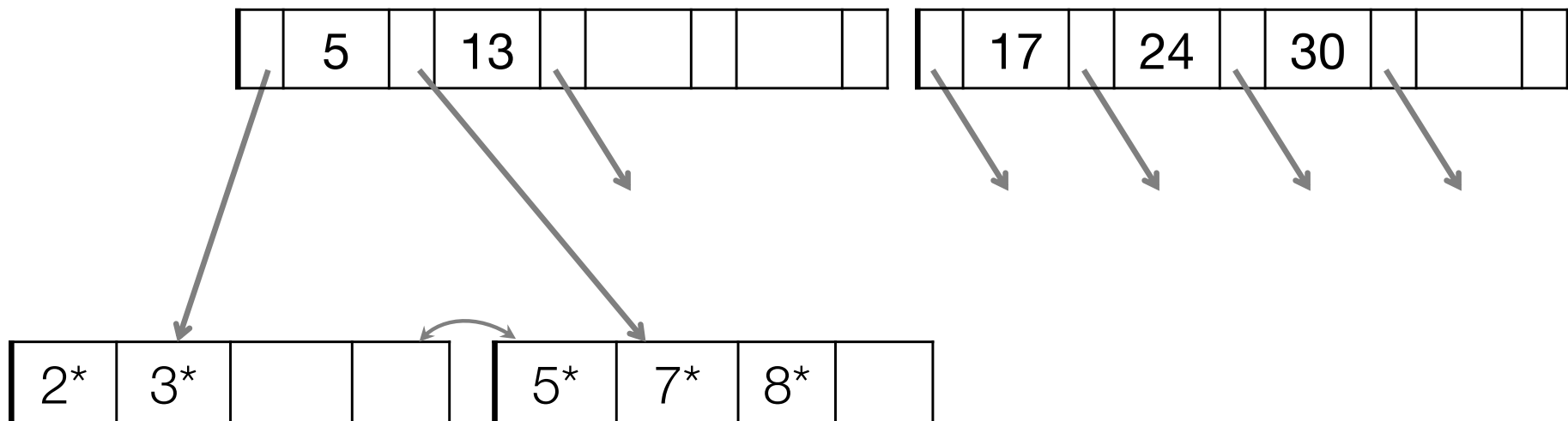
# Splitting Overfull Non-Leaf Nodes

Overfull

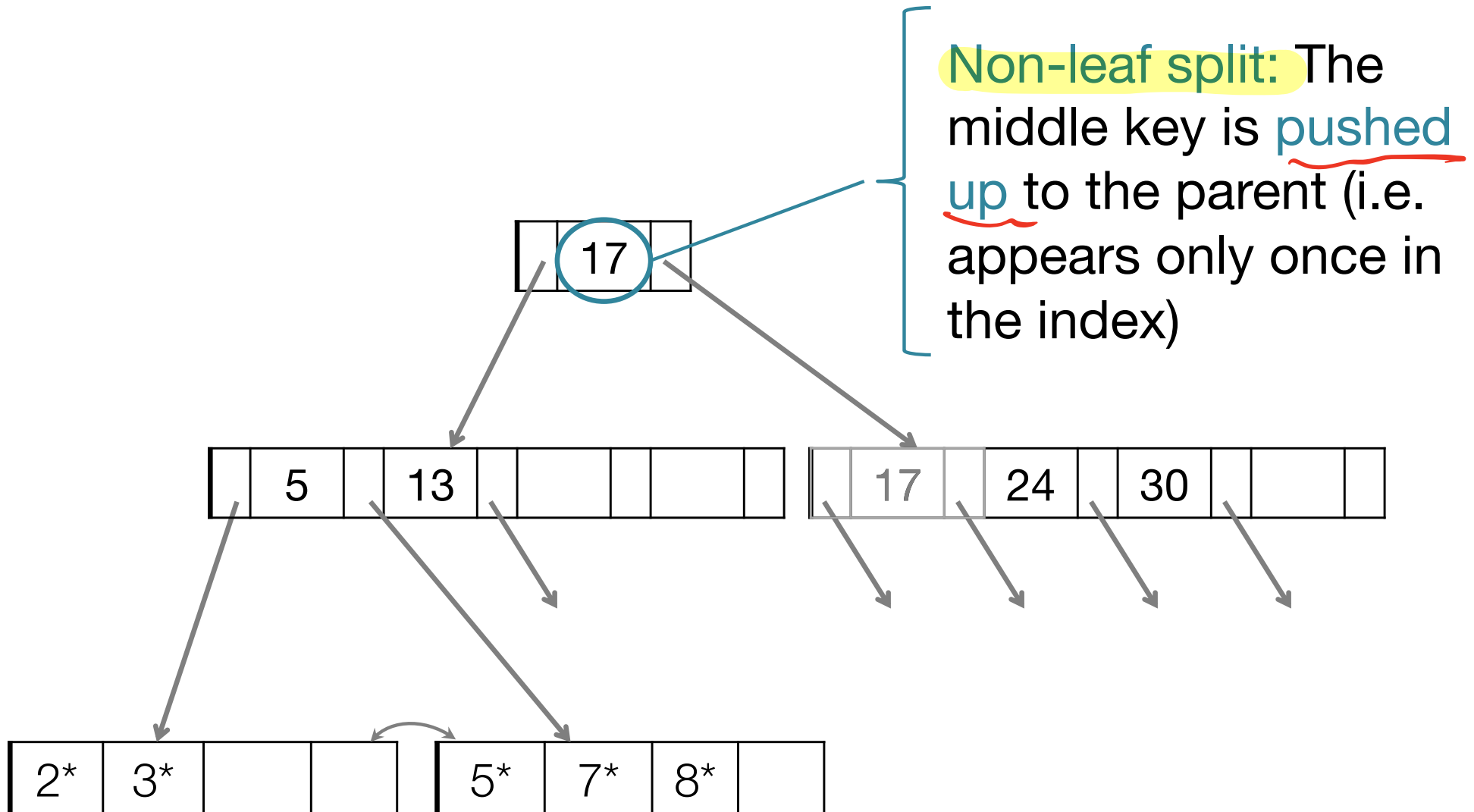


# Splitting Overfull Non-Leaf Nodes

Split

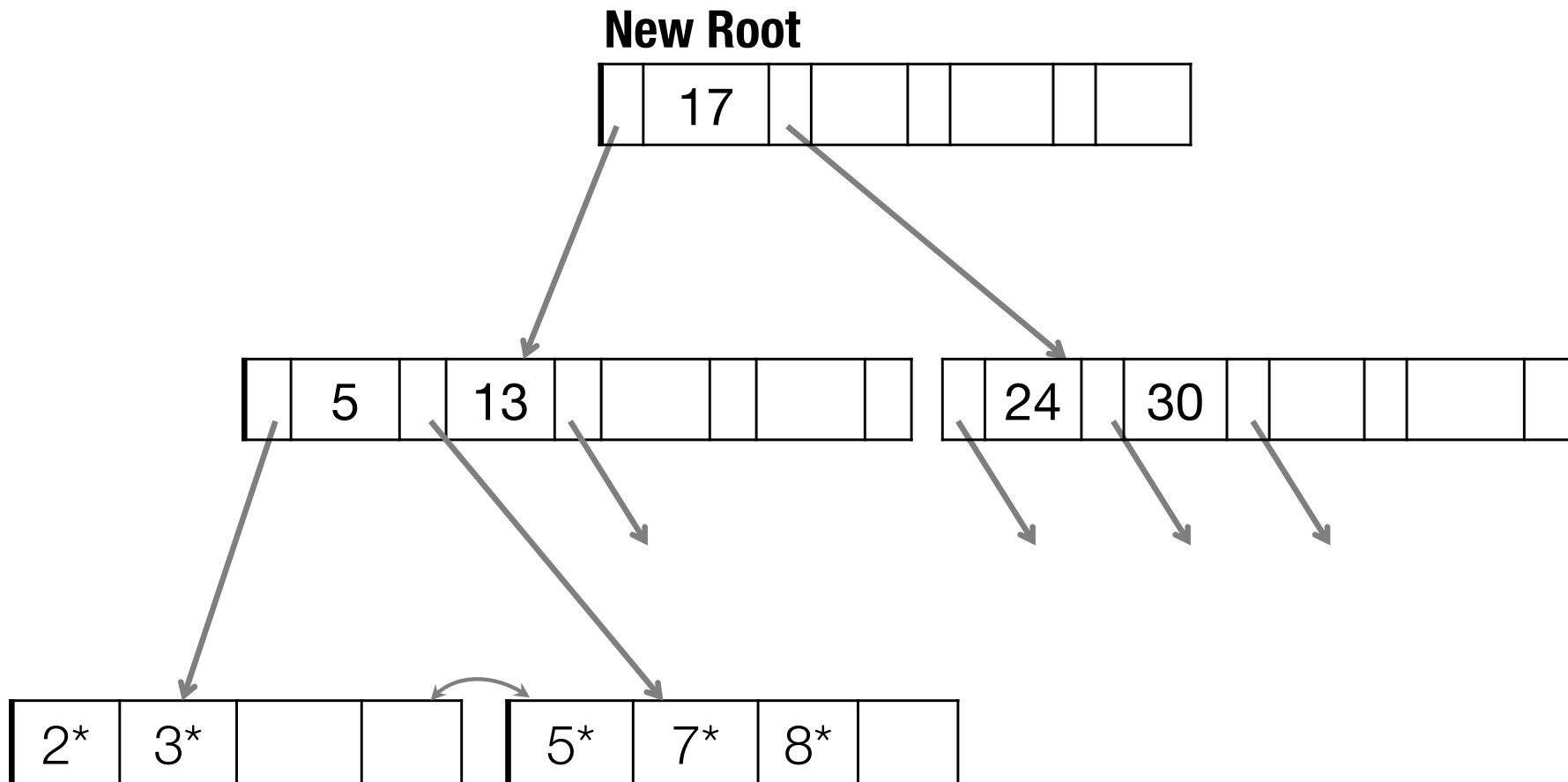


# Splitting Overfull Non-Leaf Nodes

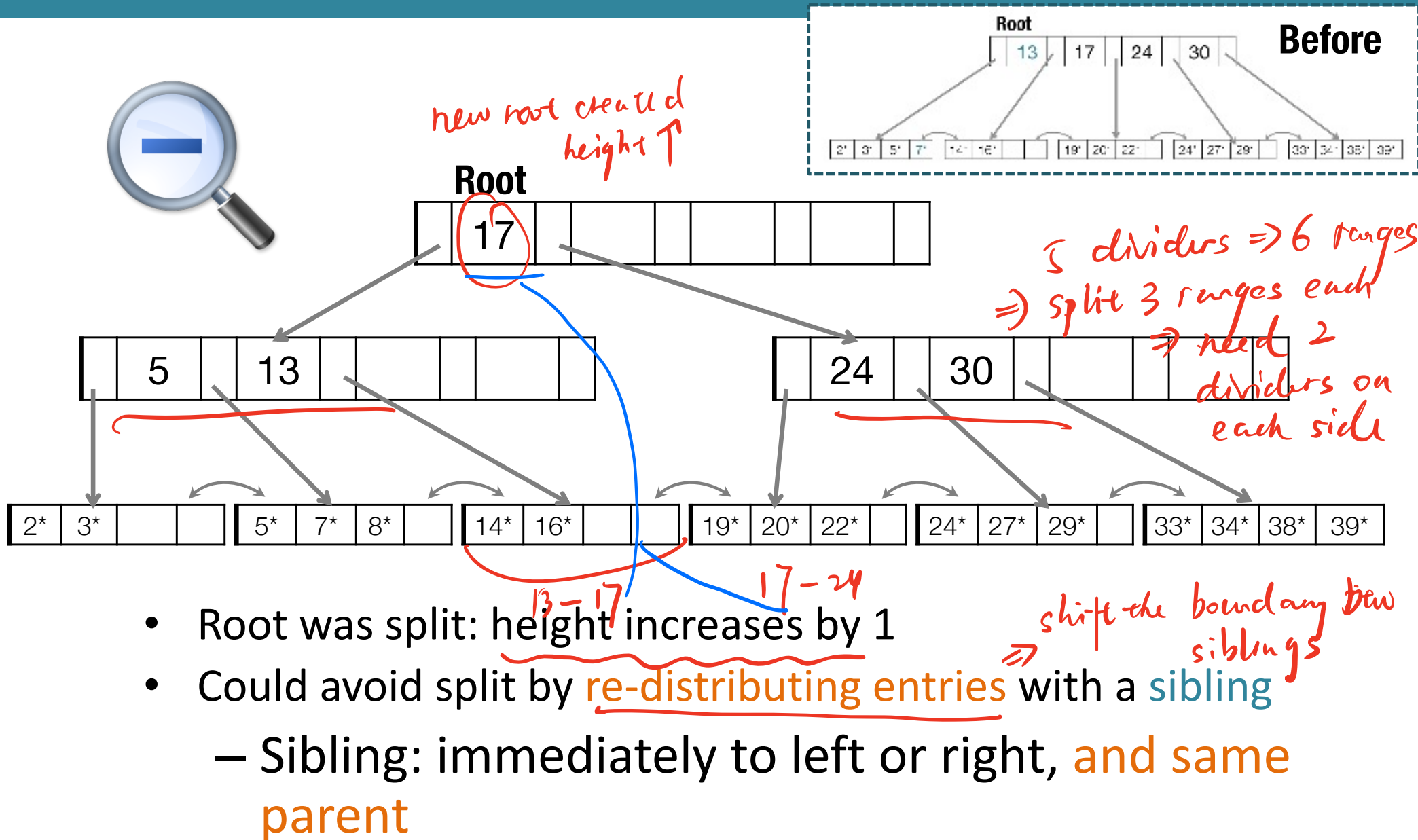


# Splitting Overfull Non-Leaf Nodes

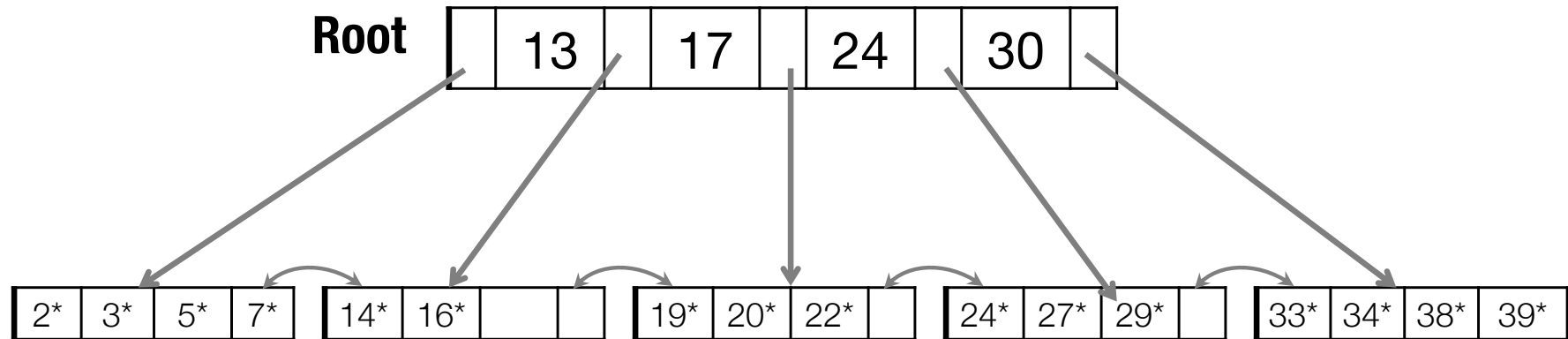
Minimum occupancy is guaranteed in both leaf and index page splits (but not in the root)



# The B<sup>+</sup>-tree After Inserting 8\*

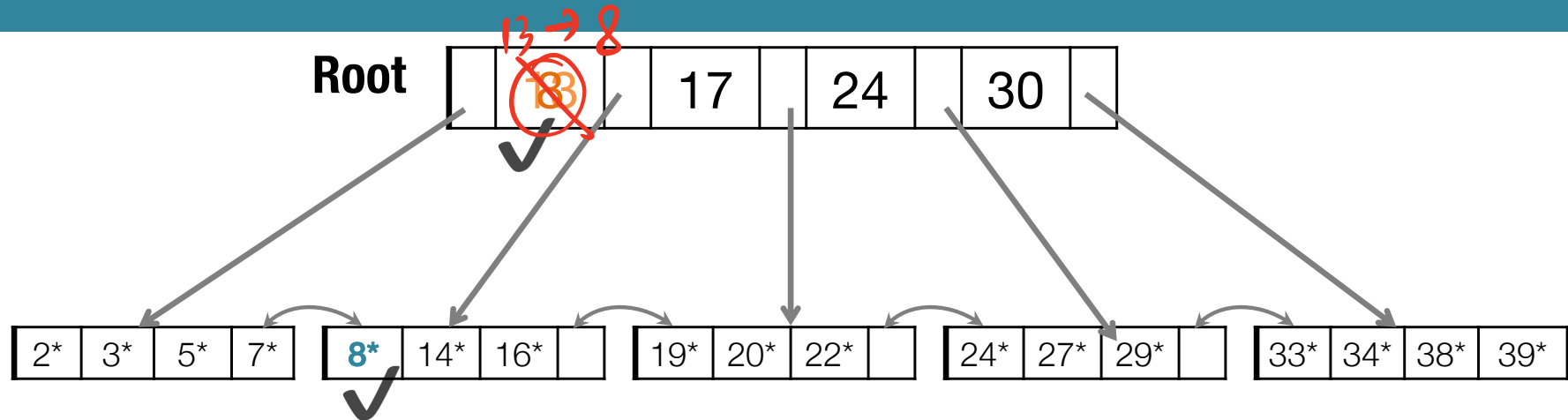


# Inserting 8\* via Entry Re-distribution with Siblings





# Inserting 8\* via Entry Re-distribution with Siblings



- Re-distributing entries with a sibling
  - Improves page occupancy, possibly reduces height
  - Usually not used for non-leaf node splits. Why?
    - Increases I/O, especially if we check both siblings
    - Can be OK to use splitting even though it propagates up the tree since that propagation would be rare and result in fewer splits on future inserts.
  - Use only for leaf level entries
    - Only have to set pointers

# Question??

Select all that are true.

The height of a B-tree

- A. Can be increased by splitting a leaf ⇒ not directly but can smg root to split
- B. Can be increased by adding a new child to a leaf X
- C. Can be increased by splitting the root ✓
- D. Cannot be changed once the tree has been built X

# B+ Tree Operations

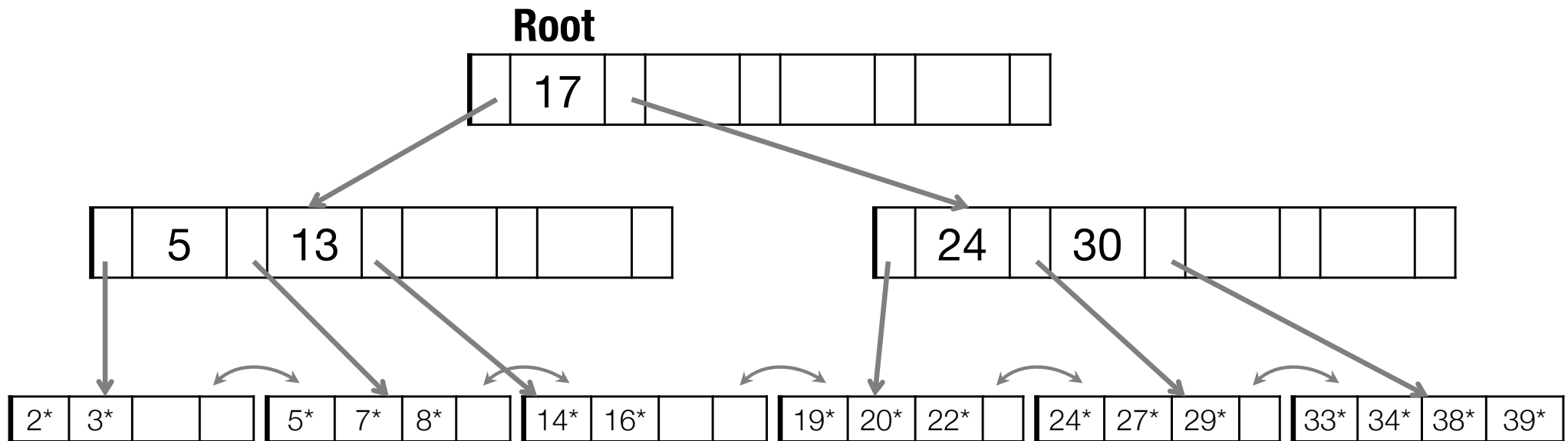
- Search
  - Equality
  - Range
- Insert data entry
- Delete data entry
- Bulk load

# B+-Tree: Deleting a Data Entry

1. Find the data entry (will always be at a leaf)
2. Delete it
3. Restore the B+ tree invariant
  - If L has **d or more** entries, done!
  - If L has only **d-1** entries,
    - ① Try to **re-distribute**, borrowing from a **sibling with more than d entries** (sibling is an adjacent node with same parent as L)
    - ② If re-distribution fails, **merge** L and sibling (Should always work!)
4. On merge, delete relevant entry in parent  
Merge could propagate to root, decreasing height.

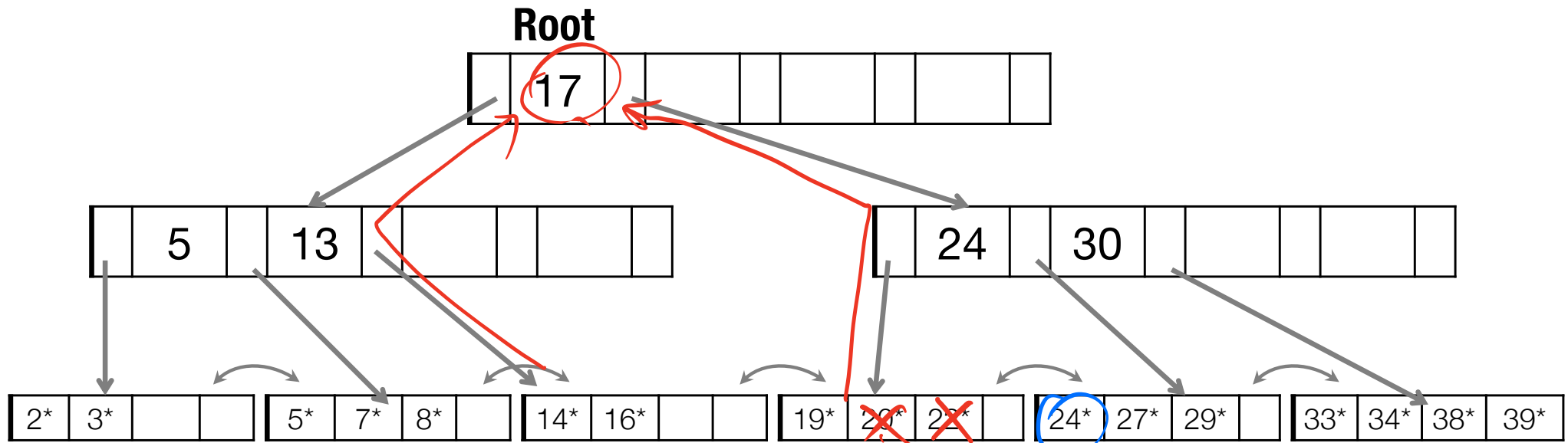


# Example Tree



- **Task:** Delete 22, 20, 24 in sequence
- Deleting 22 is easy. Invariant maintained.
- Deleting 20 is harder. Node would become less than half full.

# Deleting 22\* and 20\*



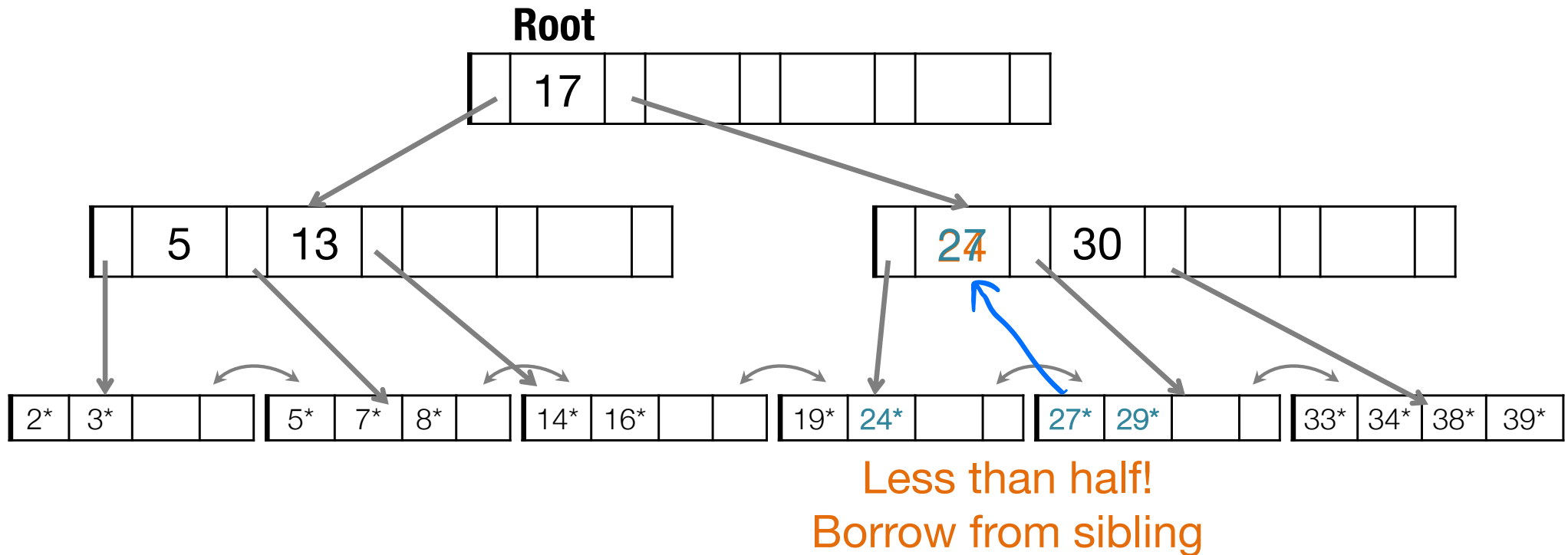
Less than half!  
Borrow from sibling

if borrow from left  
 → find common ancestor

minimum 2 not satisfied

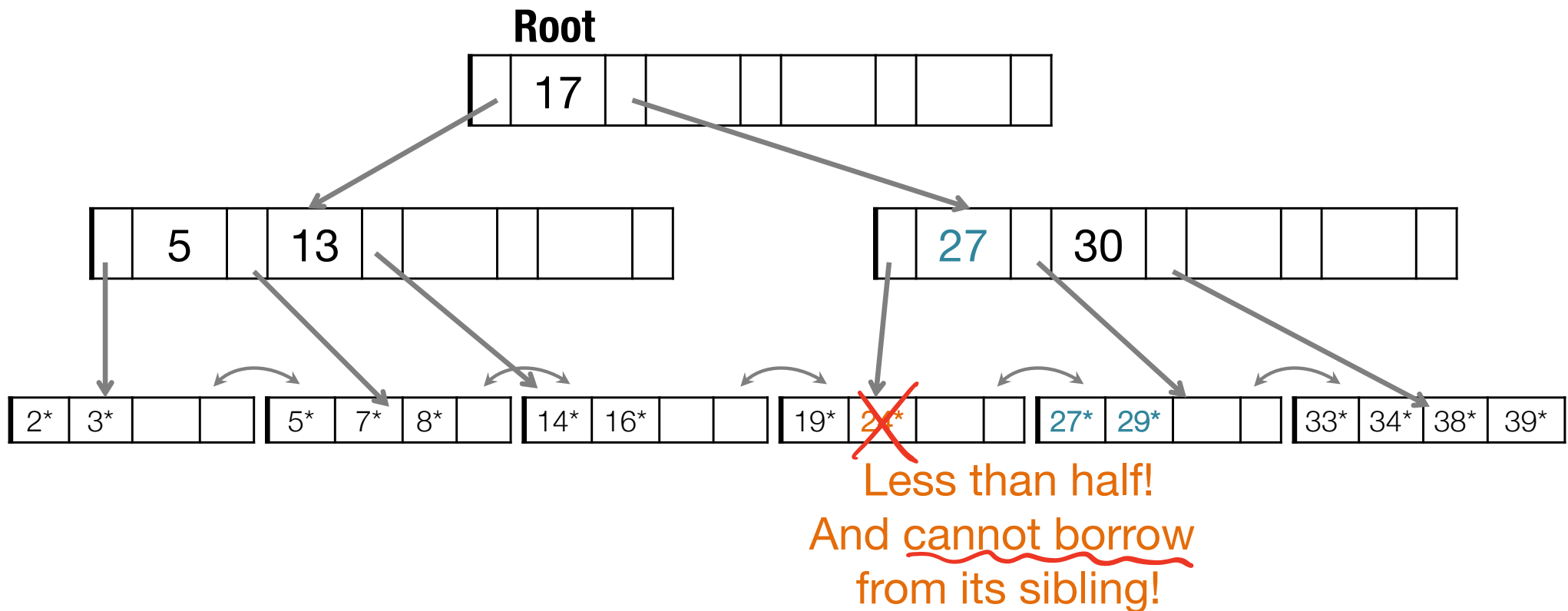
both side is OK  
 ① borrow from  
 whoever has more  
 ② borrow from  
 arbitrary side to  
 minimize 2/0

# Deleting 22\* and 20\*



Deleting 20\* is done with re-distribution.  
Notice how middle key is copied up.

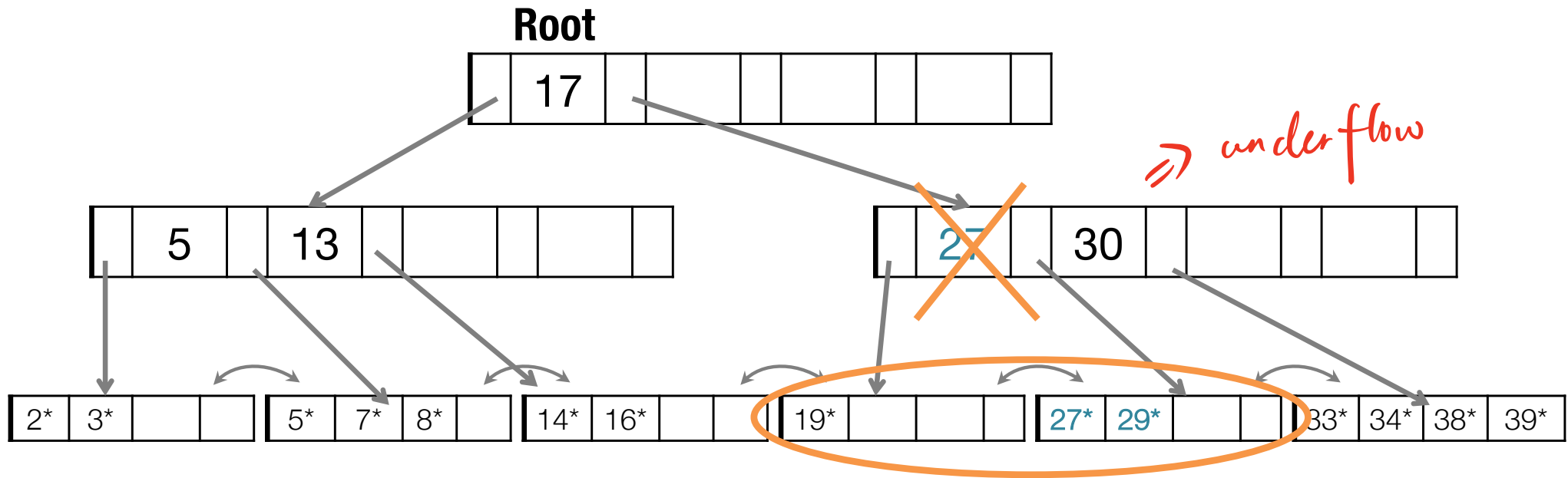
# ... And then Deleting 24\*



- Must merge  $\Rightarrow$  with sibling that you share parent with
- In the non-leaf node,  
toss the index entry with key value = 27



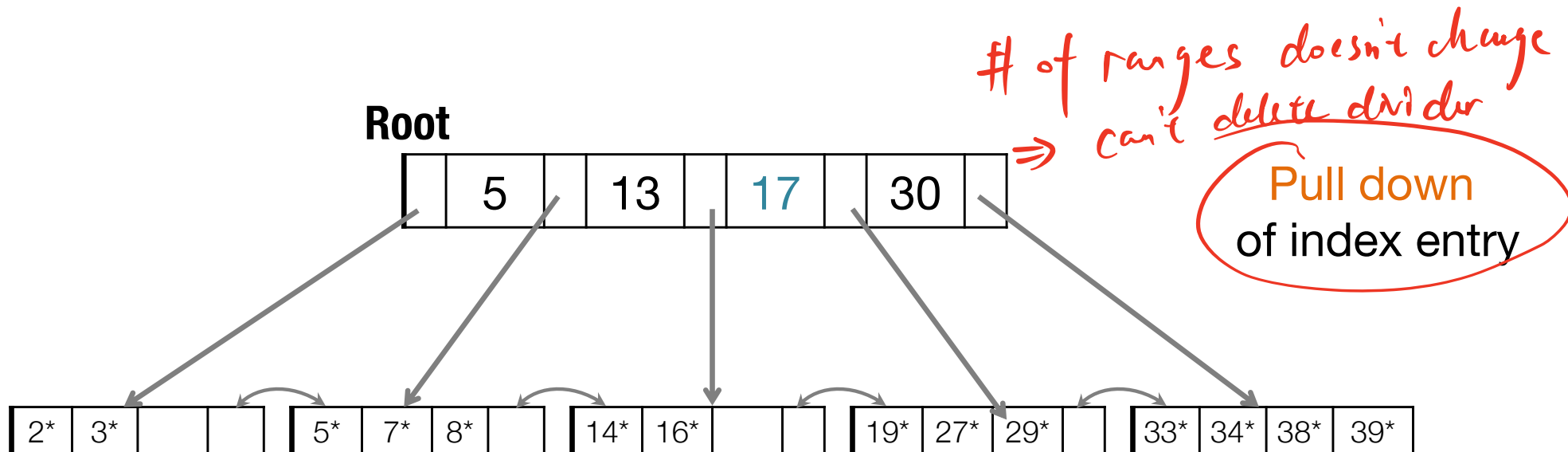
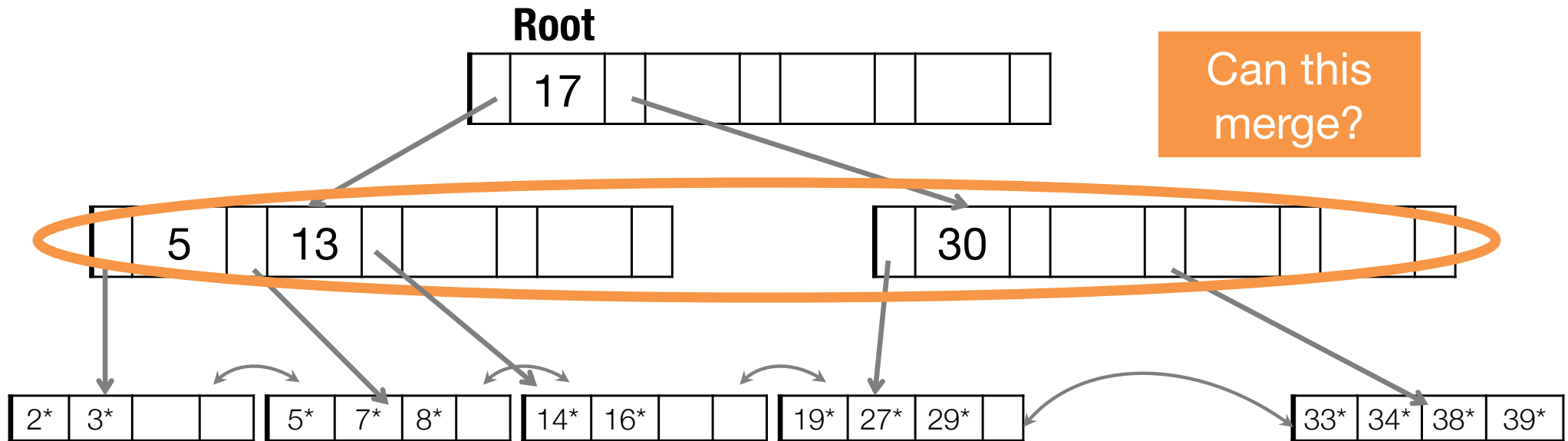
# ... And then Deleting 24\*



- Must merge
- In the non-leaf node, **toss the** index entry with key value = 27

two ranges into one  
 $\Rightarrow$  range divider btw  
has to go away

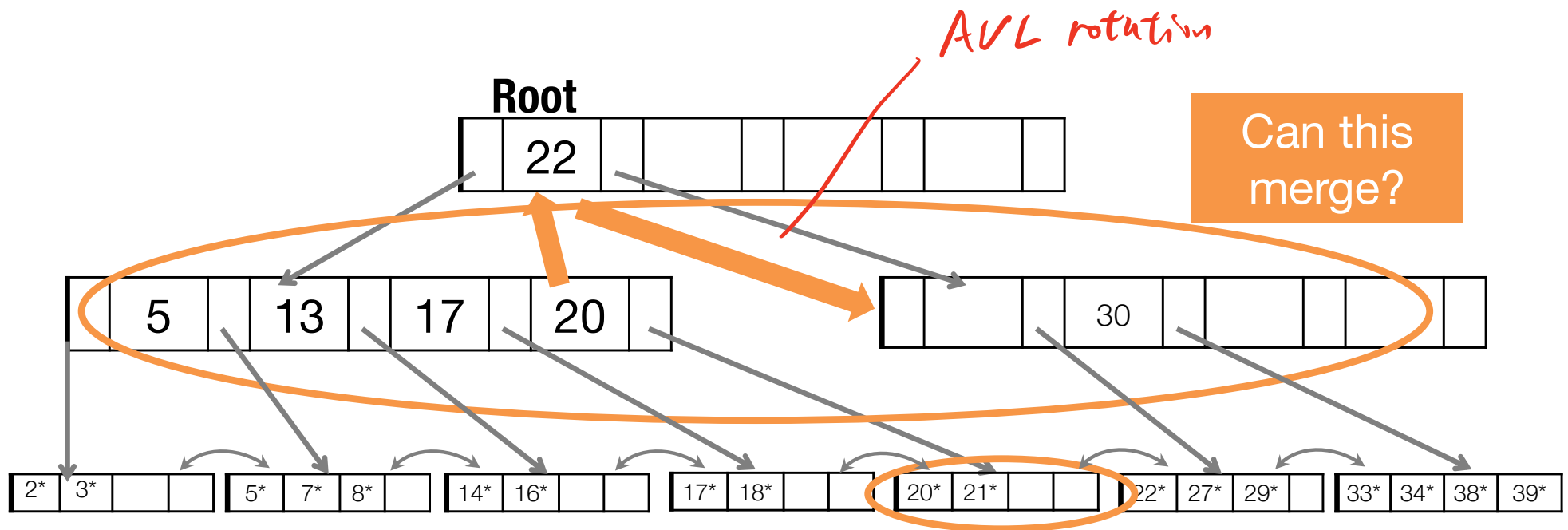
# ... And then Deleting 24\*



# Another Deletion Example:

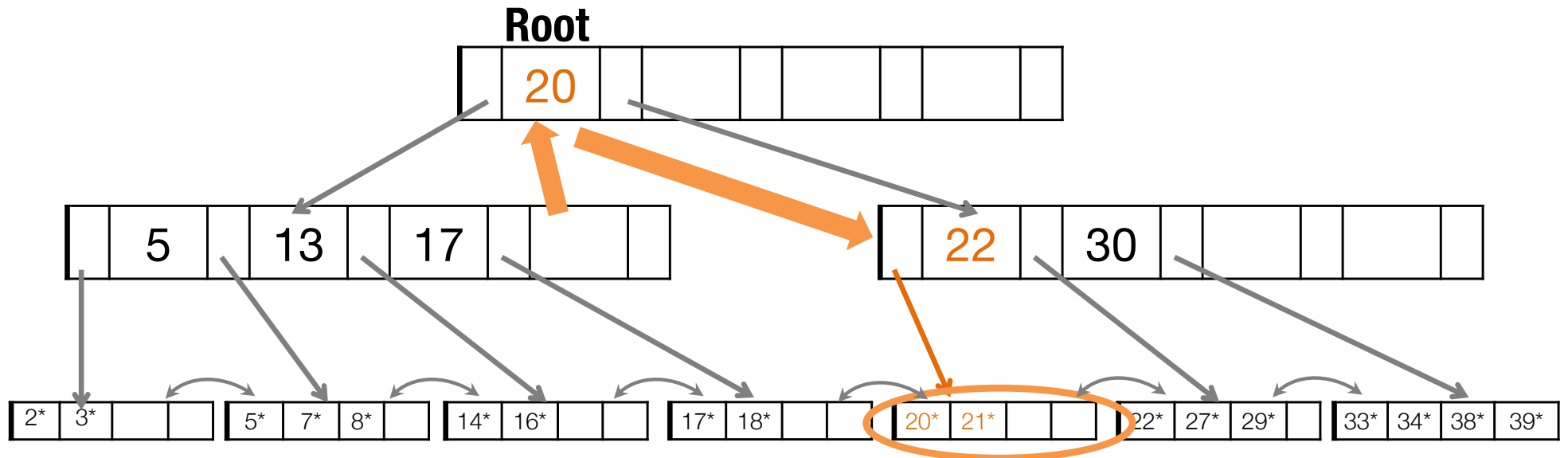
## Non-leaf Re-distribution

Can re-distribute entry from left child of root to right child.



# After Re-distribution

- Rotate through the parent node
- re-distribute index entry with key 20



# B+ Tree Deletion

- Try redistribution with **all** siblings first, then merge. Why?
  - Good chance that redistribution is possible (large fanout!)
  - Only need to propagate changes to parent node
  - Files typically grow, not shrink!

# B+ Tree Operations

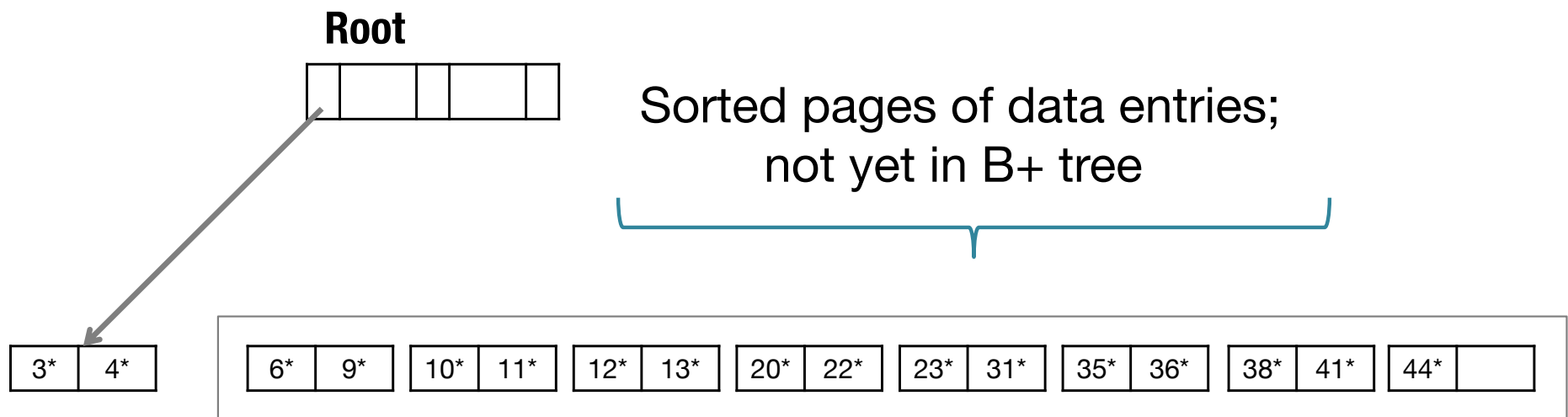
- Search
  - Equality
  - Range
- Insert data entry
- Delete data entry
- Bulk load

# Adding Entries to a B+ tree

- Option 1: multiple inserts
  - Slow. Repeated re-organization
  - Does not give sequential storage of leaves
- Option 2: Bulk Loading ✱
  - Fewer I/Os during build than one by one insert

# Bulk Loading of a B+ Tree

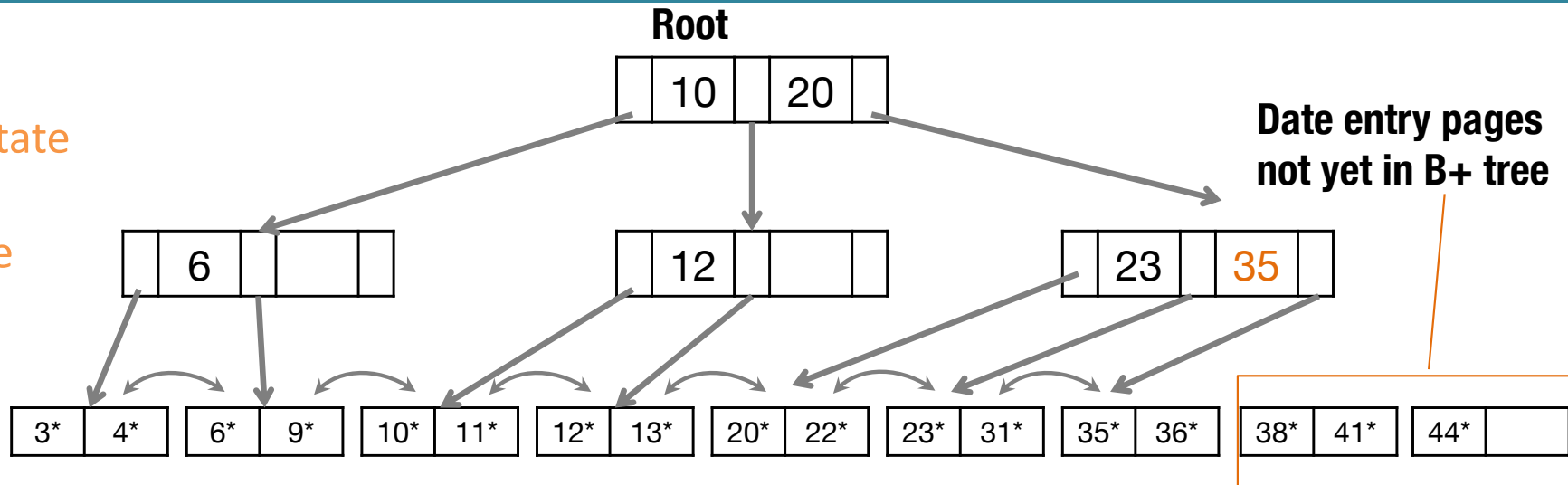
**Initialization:** Sort all data entries, insert pointer to first (leaf) page in a new (root) page. Then add one leaf page at a time, rather than one record at a time





# Bulk Loading (Contd.)

Intermediate state  
before adding  
[38\*, 41\*] page

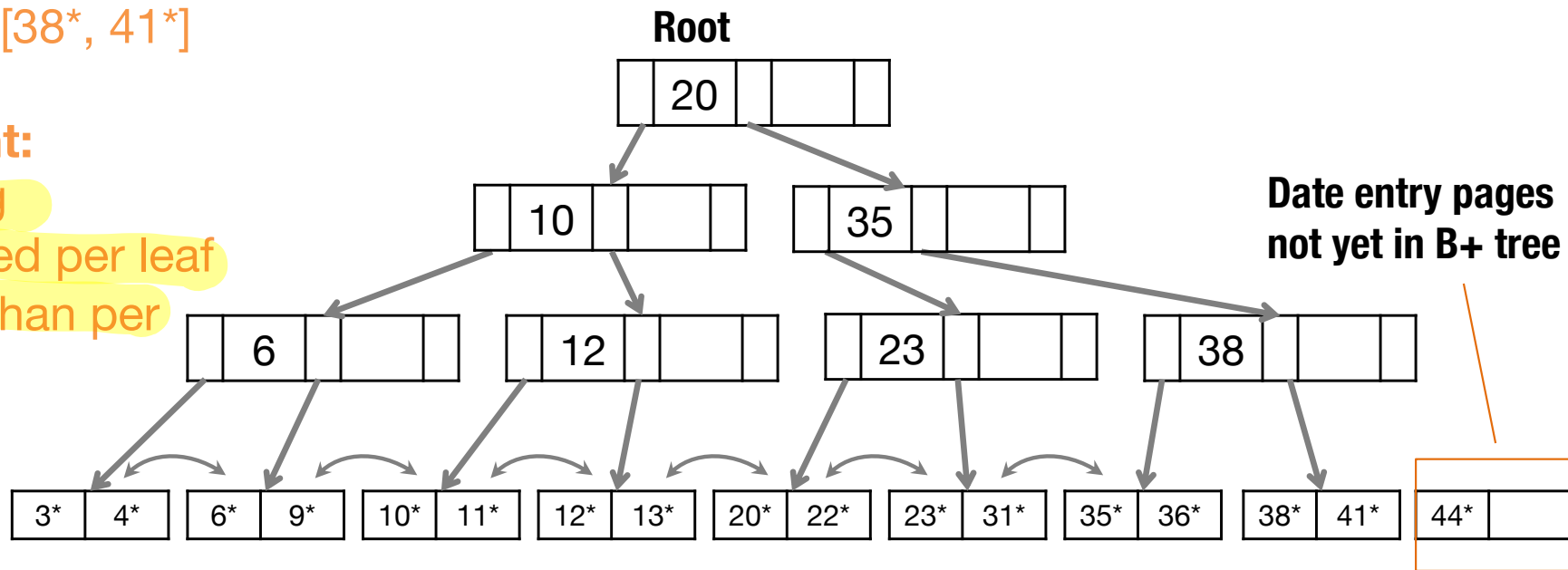


After adding [38\*, 41\*]  
to the index

**Why efficient:**

No searching

Index modified per leaf  
page rather than per  
entry



# Summary

- Tree-structured indexes are ideal for range-searches, also good for equality searches.
- B+ tree is a dynamic height-balanced index structure.
  - Insertions/deletions/search costs  $O(\log_F N)$ .
  - High fanout (F) means depth rarely more than 3 or 4. *keep top 2 in memory*  $\Rightarrow 1/2$  depth
  - Max occupancy =  $2d$ . Max fanout:  $2d+1$ .
  - Min occupancy = order =  $d$  (at least half full – exc. root)
  - Most widely used index in database management systems because of its versatility. One of the most optimized components of a DBMS.

# Suggested Exercises + Readings

Suggested Exercises: 10.1, 10.5, 10.7