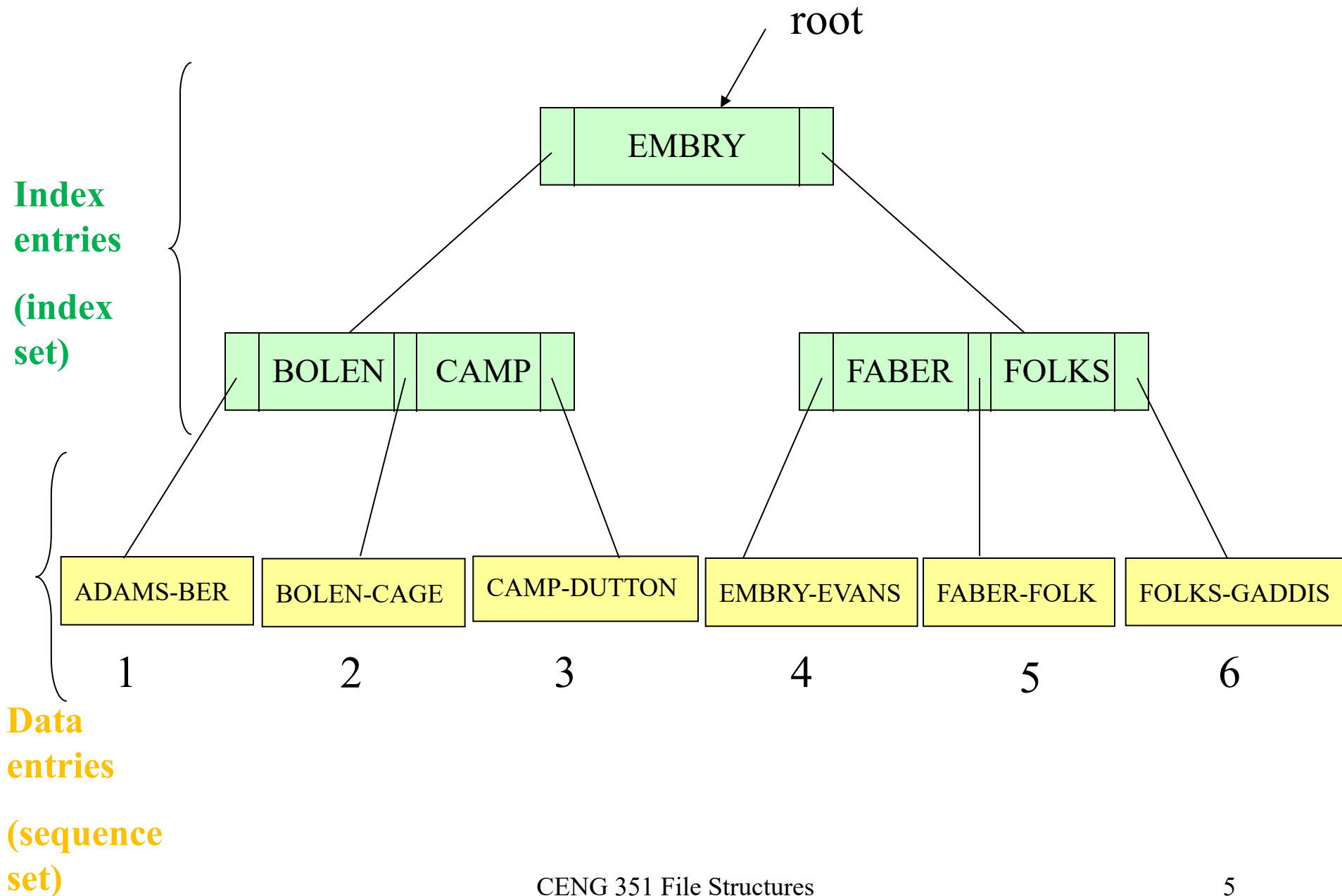# Indexing

# Part 2

# Tree indexes

- If index doesn't fit in memory:
  - Divide the index structure into blocks,
  - Organize these blocks similarly building a tree structure.

- Tree indexes:
  - B Trees
  - B+ Trees
  - Simple prefix B+ Trees
  - …

# B+ Trees

- B-tree is one of the most important data structures in computer science.

- What does B stand for? (Not binary!)

- B-tree is a **multiway search** tree.

- Several versions of B-trees have been proposed, but only B+ Trees have been used with large files.

- A B+tree is a B-tree in which data records are in leaf nodes, and faster sequential access is possible.
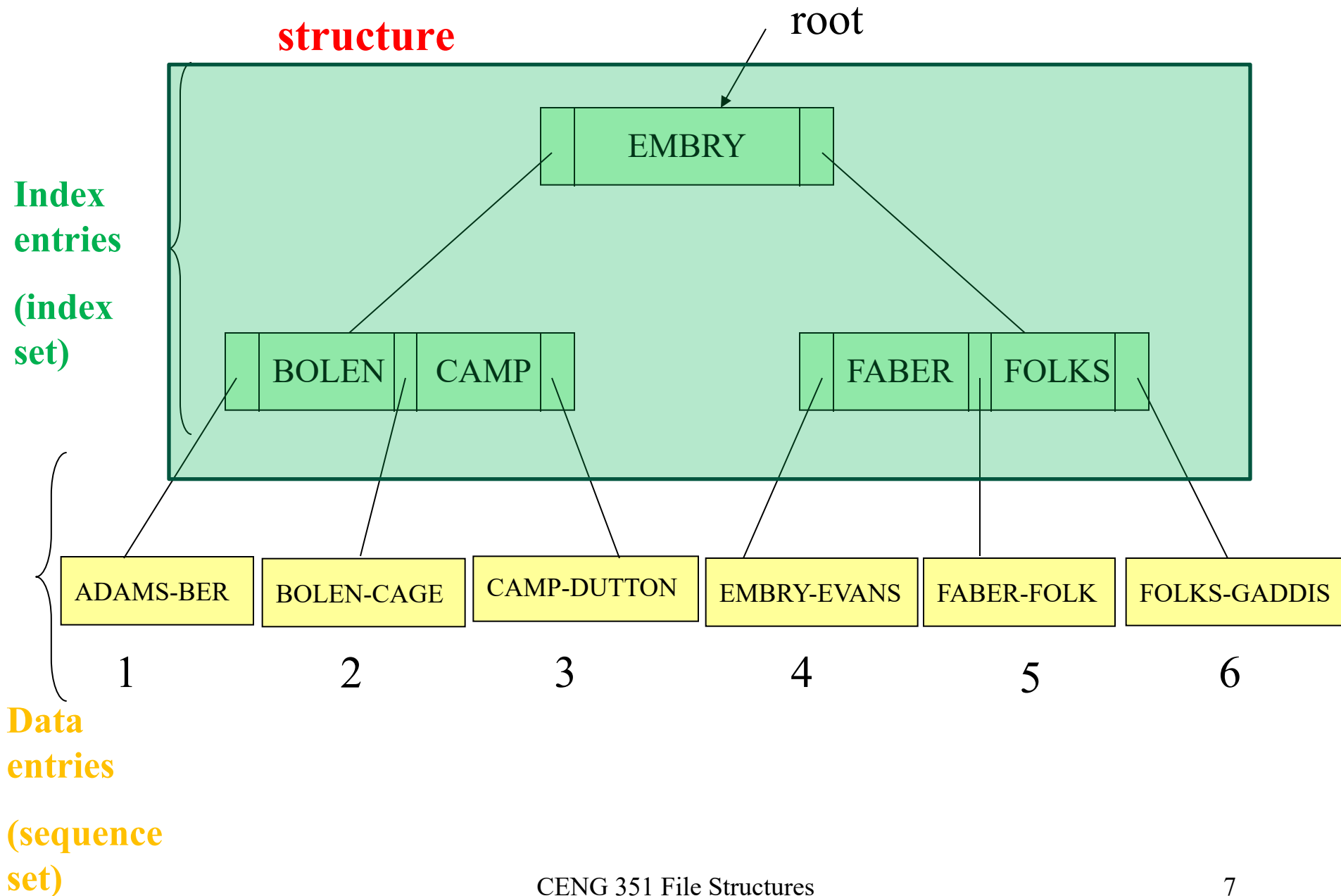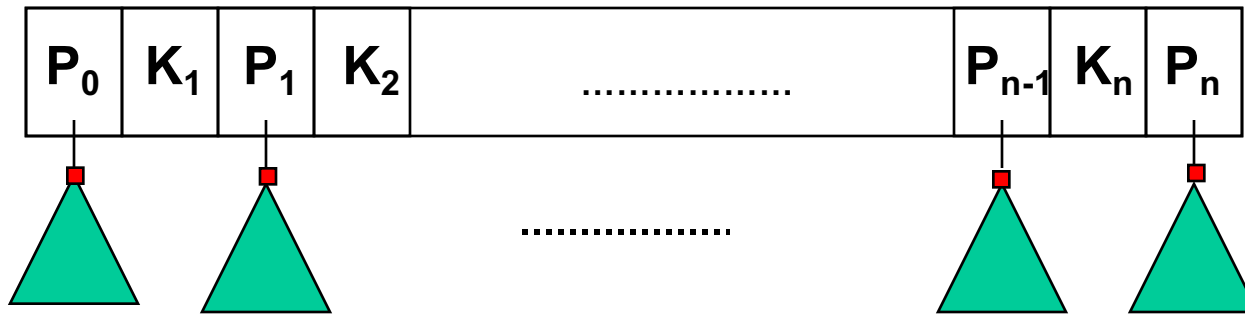
root

**Index entries (index set)**

EMBRY

BOLEN | CAMP

FABER | FOLKS

**Data entries (sequence set)**

| ADAMS-BER | BOLEN-CAGE | CAMP-DUTTON | EMBRY-EVANS | FABER-FOLK | FOLKS-GADDIS |

1     2     3     4     5     6

# Formal definition of B+ Tree Properties

- Properties of a B+ Tree of order **d**:
    - All internal nodes (except root) have **at least d keys** and **at most 2d** keys.
    - Root can have at least **1 key** and **at most 2d** keys.
    - An internal node with **n** keys has **n+1** children
    - The root has at least 2 children unless it's a leaf.
    - All leaves are on the same level (balanced tree).

# B+ tree: Internal/root node structure



**Index entries (index set)**

**Data entries (sequence set)**

root

EMBRY

BOLEN | CAMP

FABER | FOLKS

| ADAMS-BER | BOLEN-CAGE | CAMP-DUTTON | EMBRY-EVANS | FABER-FOLK | FOLKS-GADDIS |
| 1 | 2 | 3 | 4 | 5 | 6 |

# B+ tree: Internal/root node structure

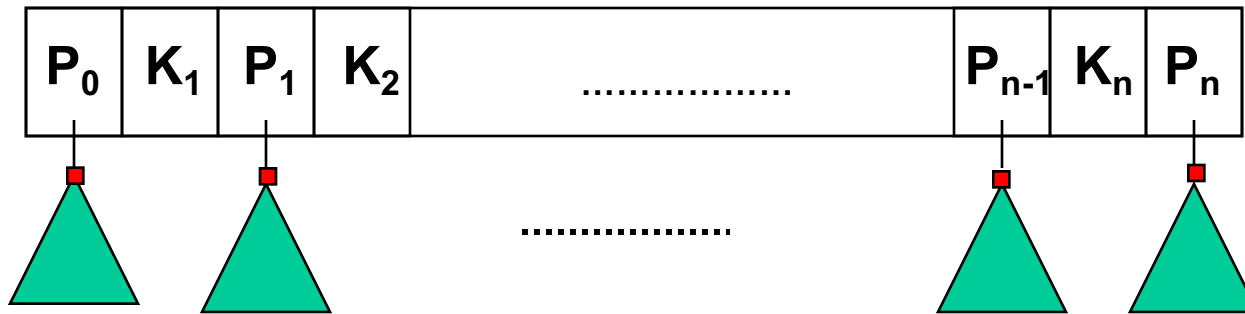| $P_0$ | $K_1$ | $P_1$ | $K_2$ | ................ | $P_{n-1}$ | $K_n$ | $P_n$ |
|---|---|---|---|---|---|---|---|

................

Each $P_i$ is a pointer to a child node; each $K_i$ is a search key value
# of search key values = n,  # of pointers = n+1

In a B+ Tree of order d:
- All internal nodes (except root) have **at least d keys** and **at most 2d** keys ( $d \le n \le 2d$).
- Root can have at least **1 key** and **at most 2d** keys. ($1 \le n \le 2d$).
- An internal node with **n** keys has **n+1** children.
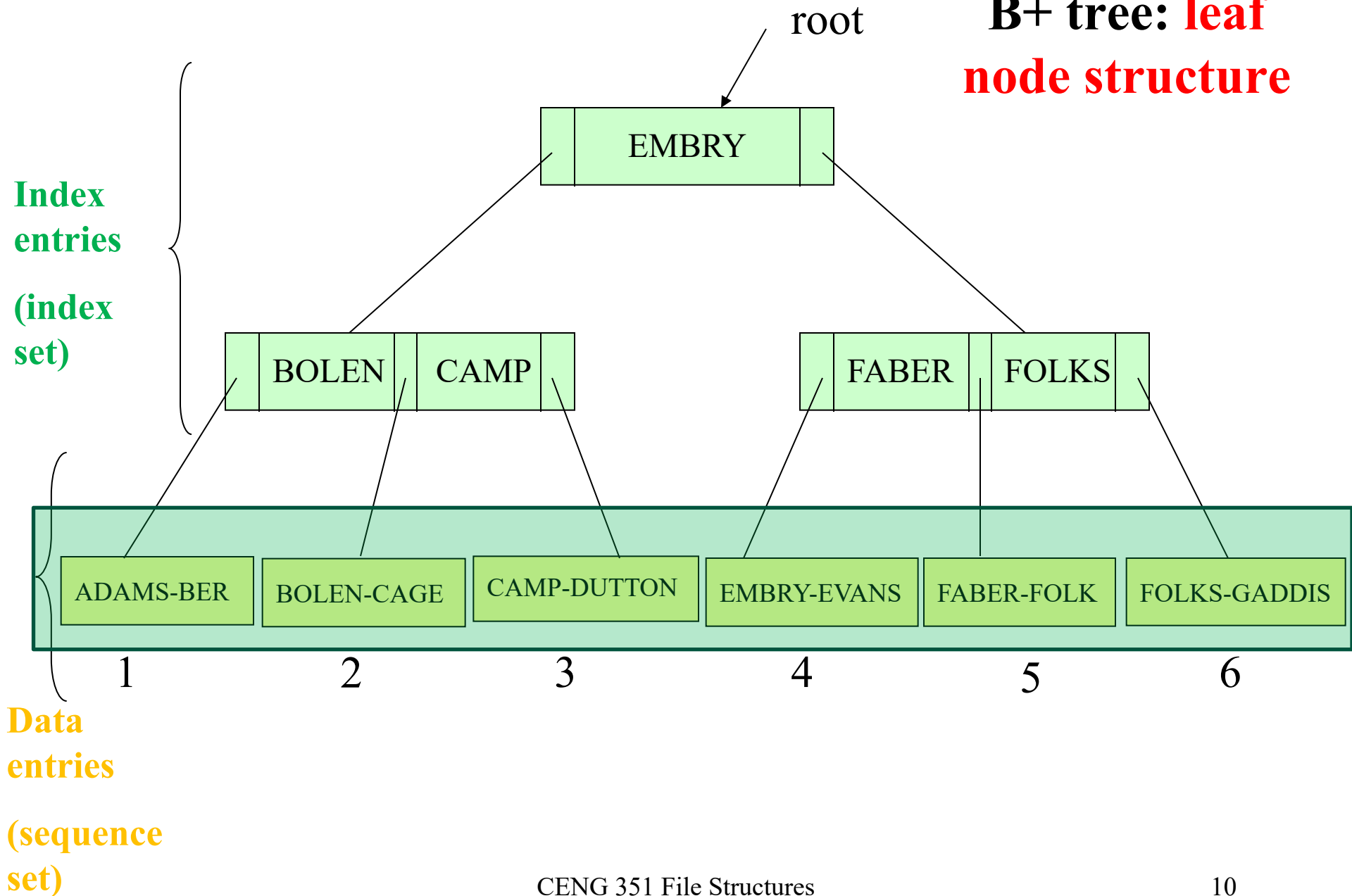- The root has at least 2 children unless it's a leaf.

8

# B+ tree: Internal/root node structure

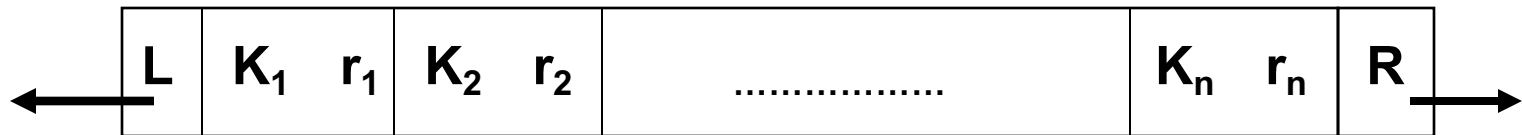| P$_0$ | K$_1$ | P$_1$ | K$_2$ | ................. | P$_{n-1}$ | K$_n$ | P$_n$ |
|---|---|---|---|---|---|---|---|

.................

Each P$_i$ is a pointer to a child node; each K$_i$ is a search key value
# of search key values = n,         # of pointers = n+1

- Requirements:
- $K_1 < K_2 < \ldots < K_n$
- For any search key value K in the subtree pointed by Pi,
  If $P_i = P_0$, we require $K < K_1$
  If $P_i = P_n$, $K_n <= K$
  If $P_i = P_1, \ldots, P_{n-1}$, $K_i <= K < K_{i+1}$
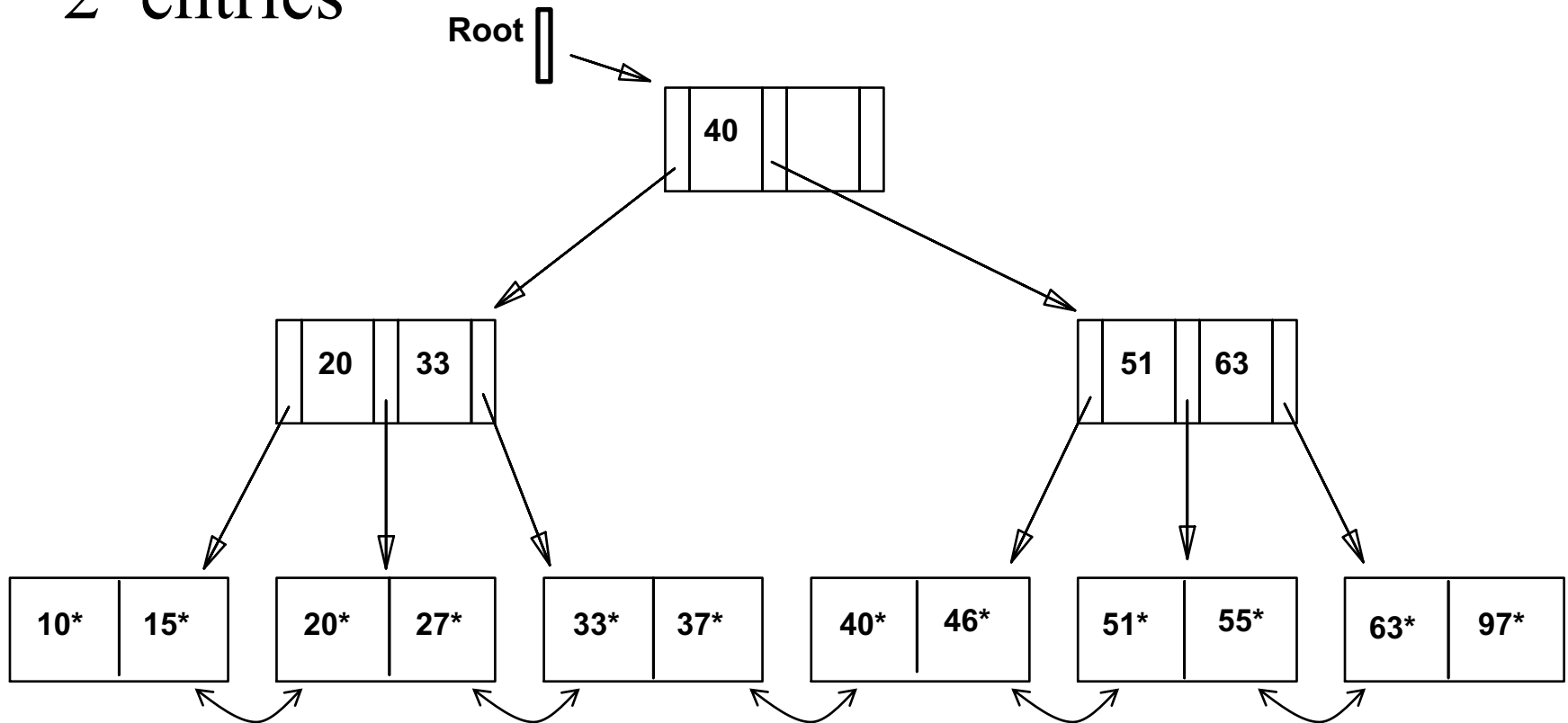
root

**B+ tree: leaf node structure**

**Index entries**

**(index set)**

EMBRY

BOLEN | CAMP

FABER | FOLKS

ADAMS-BER | BOLEN-CAGE | CAMP-DUTTON | EMBRY-EVANS | FABER-FOLK | FOLKS-GADDIS

1    2    3    4    5    6

**Data entries**

**(sequence set)**

# B+ tree: leaf node structure

| L | $K_1$  $r_1$ | $K_2$  $r_2$ | ................. | $K_n$  $r_n$ | R |
|---|---|---|---|---|---|

- Pointer L points to the left neighbor; R points to the right neighbor (**doubly linked list**)

- $K_1 < K_2 < \ldots < K_n$

- $d \leq n \leq 2d$ (d is the order of this B+ tree)
- We will use $K_i^*$ for the pair $<K_i, r_i>$ and omit L and R for simplicity
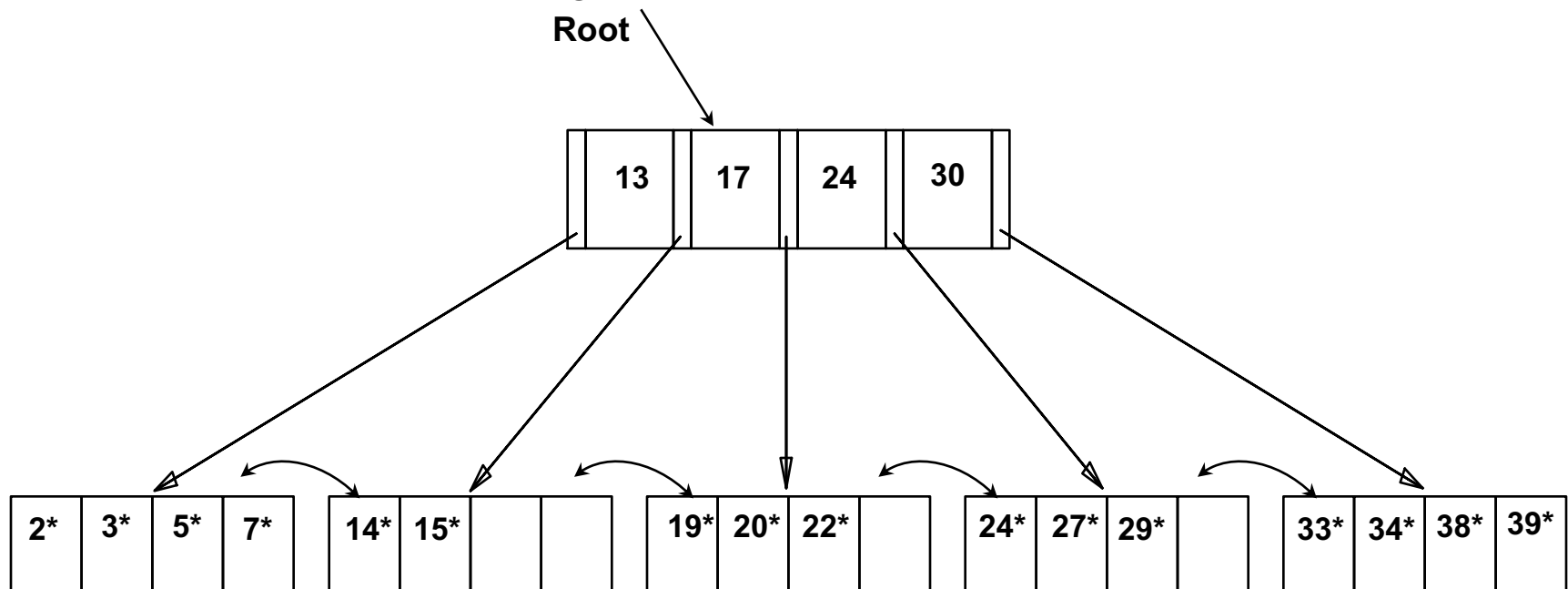- All leaves are on the same level (balanced tree).

# Example: B+ tree with order of 1

- Each node must hold at least 1 entry, and at most 2 entries

**Root**

```
              40
```

```
   20   33              51   63
```

```
10*  15*   20*  27*   33*  37*   40*  46*   51*  55*   63*  97*
```

# Example: Search in a B+ tree order 2

- Search: how to find the records with a given search key value?
  - Begin at root, and use key comparisons to go to leaf

- Examples: search for 5*, 16*, all data entries >= 24* ...
  - The last one is a range search, we need to do the sequential scan, starting from the first leaf containing a value >= 24.
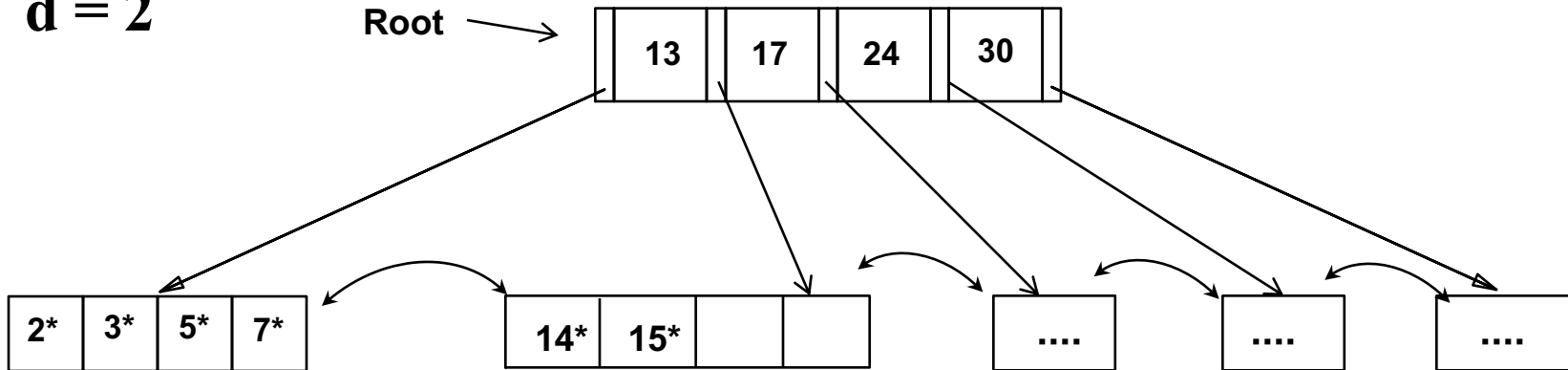
**Root**

| | 13 | | 17 | | 24 | | 30 | |

| 2* | 3* | 5* | 7* | | 14* | 15* | | | | 19* | 20* | 22* | | | 24* | 27* | 29* | | | 33* | 34* | 38* | 39* |

# How to Insert a Data Entry into a B+ Tree?

- Let's look at several examples first.

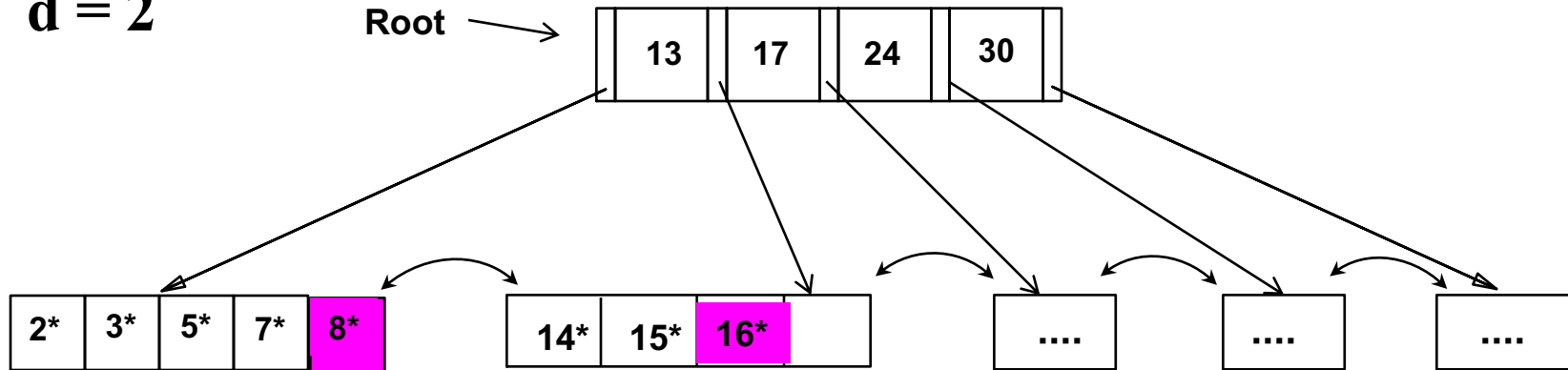# Inserting 16*, 8* into Example B+ tree

**d = 2**

**Root** →

| | 13 | 17 | 24 | 30 | |
|---|---|---|---|---|---|

| 2* | 3* | 5* | 7* |
|---|---|---|---|

| 14* | 15* | | |
|---|---|---|---|

| **....** |
|---|

| **....** |
|---|

| **....** |
|---|

Leaf nodes:
$d \leq n \leq 2d$
$2 \leq n \leq 4$

# Inserting 16*, 8* into Example B+ tree

d = 2

**Root**

| | 13 | 17 | 24 | 30 |
|---|---|---|---|---|

| 2* | 3* | 5* | 7* | 8* |
|---|---|---|---|---|

| 14* | 15* | 16* | |
|---|---|---|---|

| .... |
|---|

| .... |
|---|

| .... |
|---|

**Leaf node overflows!!!**

Leaf nodes:

$d \le n \le 2d$

$2 \le n \le 4$

# Inserting 16*, 8* into Example B+ tree

**d = 2**

**Root**

| | 13 | 17 | 24 | 30 |
|---|---|---|---|---|

| 2* | 3* | 5* | 7* | 8* |
|---|---|---|---|---|

| 14* | 15* | 16* | |
|---|---|---|---|

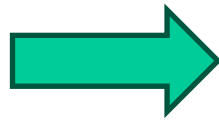| .... |
|---|

| .... |
|---|

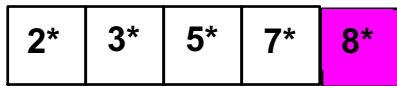| .... |
|---|

**Leaf node overflows!!!**

**When a <u>leaf node</u> overflows:**

1) Split the node

   First **d** entries stay in old node , move rest of entries  to new node

| 2* | 3* | 5* | 7* | 8* |
|---|---|---|---|---|

# Inserting 16*, 8* into Example B+ tree

**d = 2**

**Root**

| | 13 | 17 | 24 | 30 |
|---|---|---|---|---|

| 2* | 3* | 5* | 7* | 8* |
|---|---|---|---|---|

| 14* | 15* | 16* | |
|---|---|---|---|

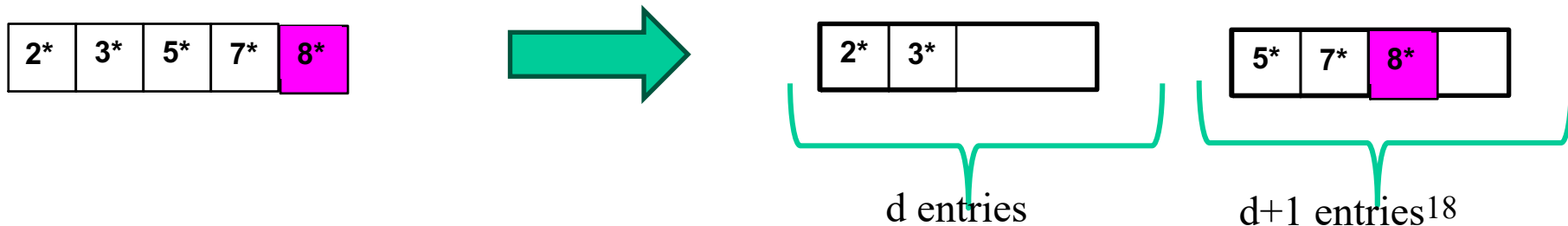| .... |
|---|

| .... |
|---|

| .... |
|---|

**Leaf node overflows!!!**

**When a leaf node overflows:**

1) Split the node

   First **d** entries stay in old node , move rest of entries  to new node

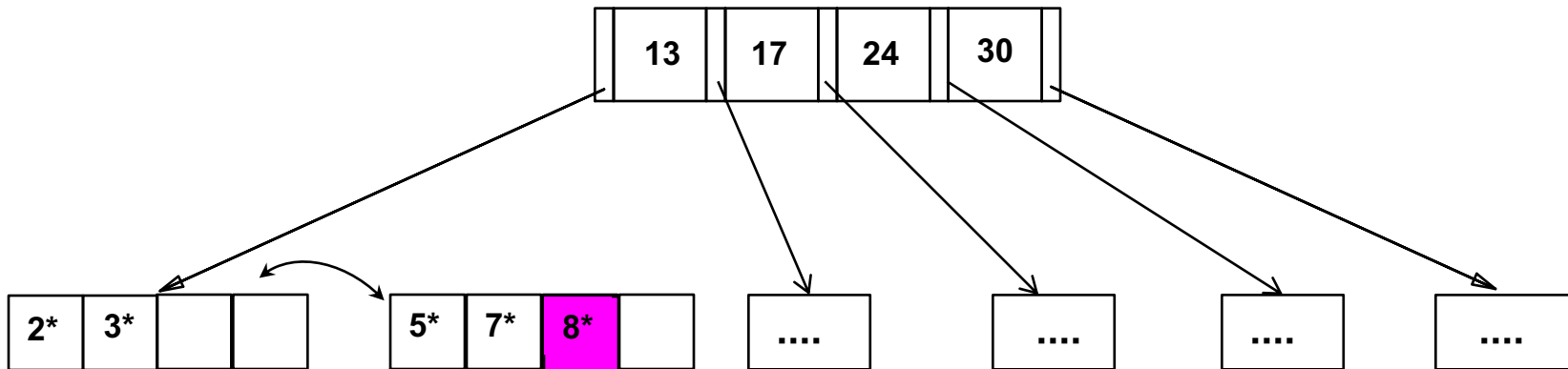| 2* | 3* | 5* | 7* | 8* |
|---|---|---|---|---|

| 2* | 3* | |
|---|---|---|

| 5* | 7* | 8* | |
|---|---|---|---|

d entries

d+1 entries[18]

# Inserting 16*, 8* into Example B+ tree

**d = 2**

**Root** →

| | 13 | 17 | 24 | 30 |
|--|----|----|----|----|

| 2* | 3* | 5* | 7* | **8*** |
|----|----|----|----|------|

| 14* | 15* | **16*** | |
|-----|-----|-------|--|

| **....** |
|------|

| **....** |
|------|

| **....** |
|------|

**Leaf node overflows!!!**          ⇓

| | 13 | 17 | 24 | 30 |
|--|----|----|----|----|

| 2* | 3* | | |
|----|----|--|--|

| 5* | 7* | **8*** | |
|----|----|------|--|

| **....** |
|------|

| **....** |
|------|

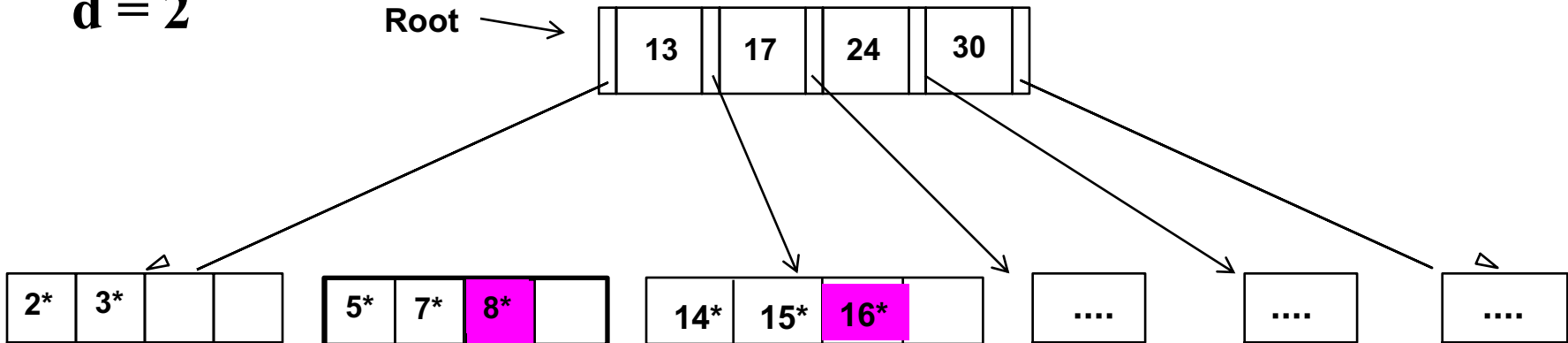| **....** |
|------|

| **....** |
|------|

One new child (leaf node)
generated; must add **one more
pointer** to its parent, thus **one
more key value** as well**.**

19

# Inserting 16*, 8* into Example B+ tree

**d = 2**



**Leaf node overflows!!!**
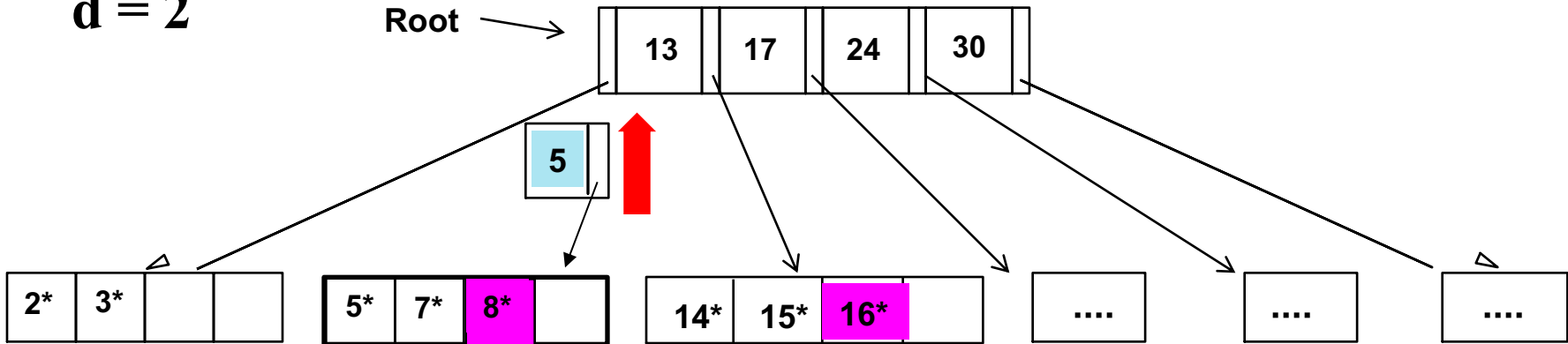
**When a <u>leaf node</u> overflows:**

1) Split the node
   First **d** entries stay, move rest to new node

2) We need a pointer to the new block for the search: **COPY UP** the *middle key* as the search key. Also, add pointer to the new block

# Inserting 16*, 8* into Example B+ tree

d = 2

Root

| | 13 | 17 | 24 | 30 |

5

| 2* | 3* | | |

| 5* | 7* | 8* | |

| 14* | 15* | 16* | |

| .... |

| .... |

| .... |

**Leaf node overflows!!!**

**When a <u>leaf node</u> overflows:**
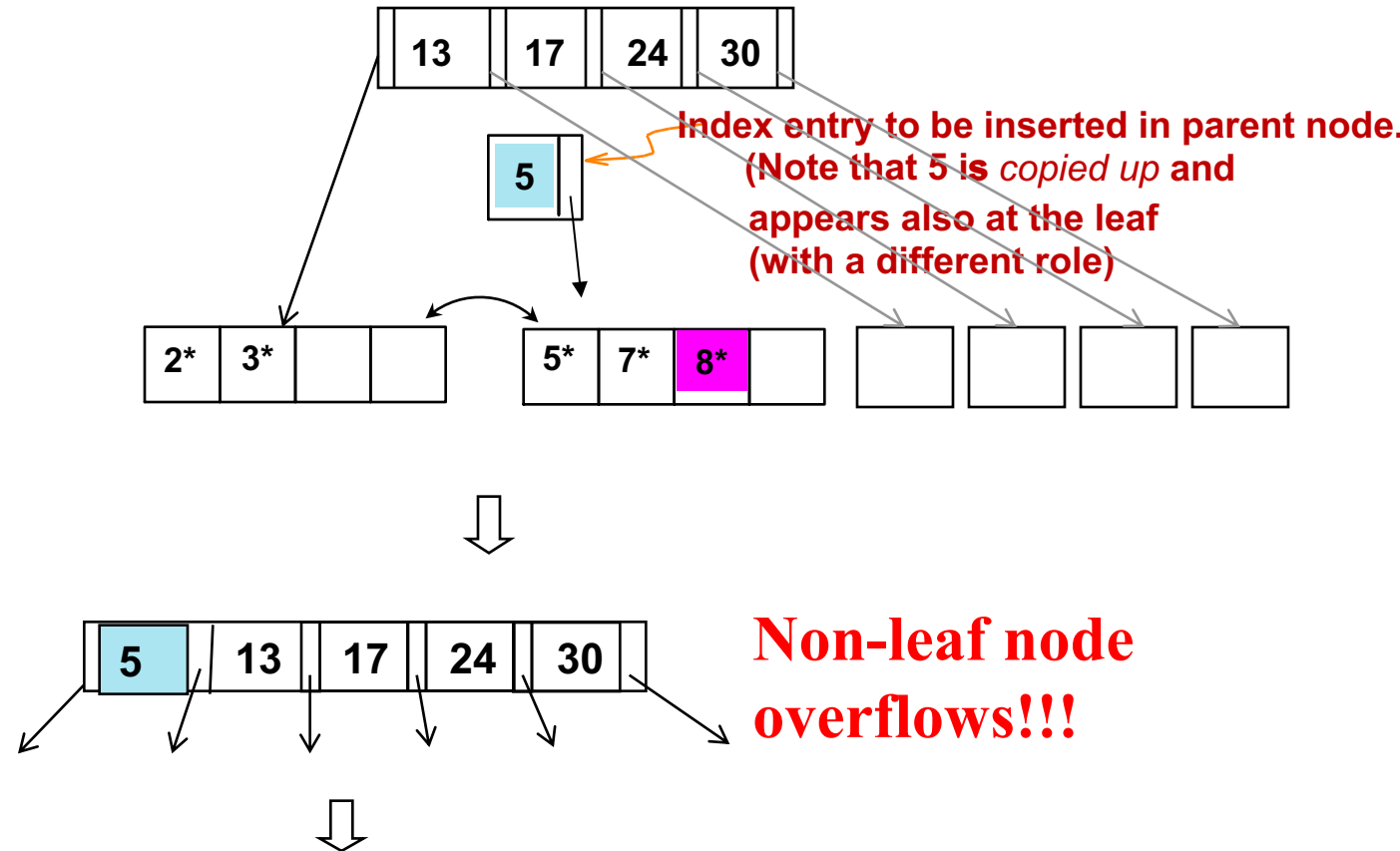
1) Split the node
   First **d** entries stay, move rest to new node

2) **COPY UP** the *middle key* as the search key. Also, add pointer to the new block
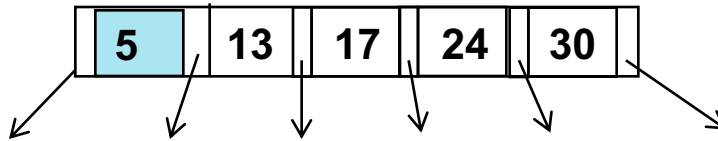
# Inserting 8* (cont.)

**d = 2**

- Copy up the middle value (leaf split)



13 | 17 | 24 | 30

5

Index entry to be inserted in parent node.
(Note that 5 is *copied up* and appears also at the leaf (with a different role)

2* | 3*

5* | 7* | 8*

5 | 13 | 17 | 24 | 30

**Non-leaf node overflows!!!**

# Inserting 8* (cont.)

**d** = 2

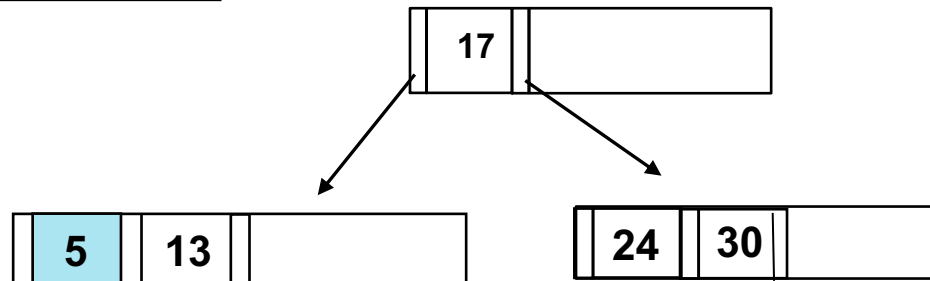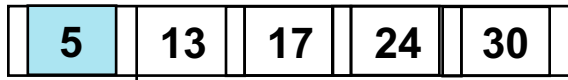| 5 | 13 | 17 | 24 | 30 |
|---|----|----|----|----|

**Non-leaf node overflows!!!**

**When a non-leaf node overflows:**

1) Split the node

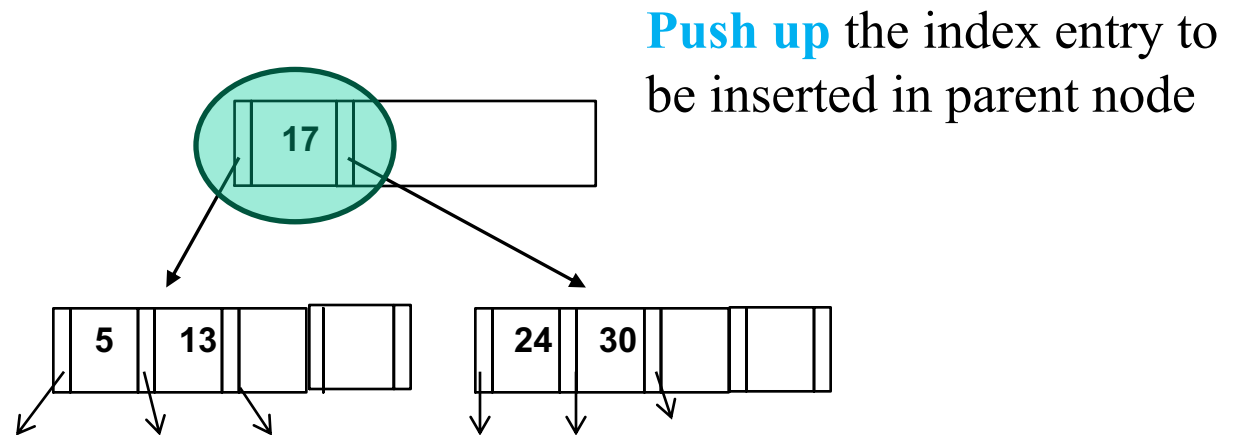| First **d** keys (and d+1 pointers) stay in old node) | Last **d** keys (and d+1 pointers) move to new node |
|---|---|

2) **PUSH UP** middle key (17) and (pointers to the blocks)!

| 5 | 13 | 17 | 24 | 30 |
|---|----|----|----|----|

➡

| 17 |
|----|

| 5 | 13 |
|---|----|

| 24 | 30 |
|----|----|

# Insertion into B+ tree (cont.)

Recall: Root can have at least **1 key** and **at most 2d** keys.
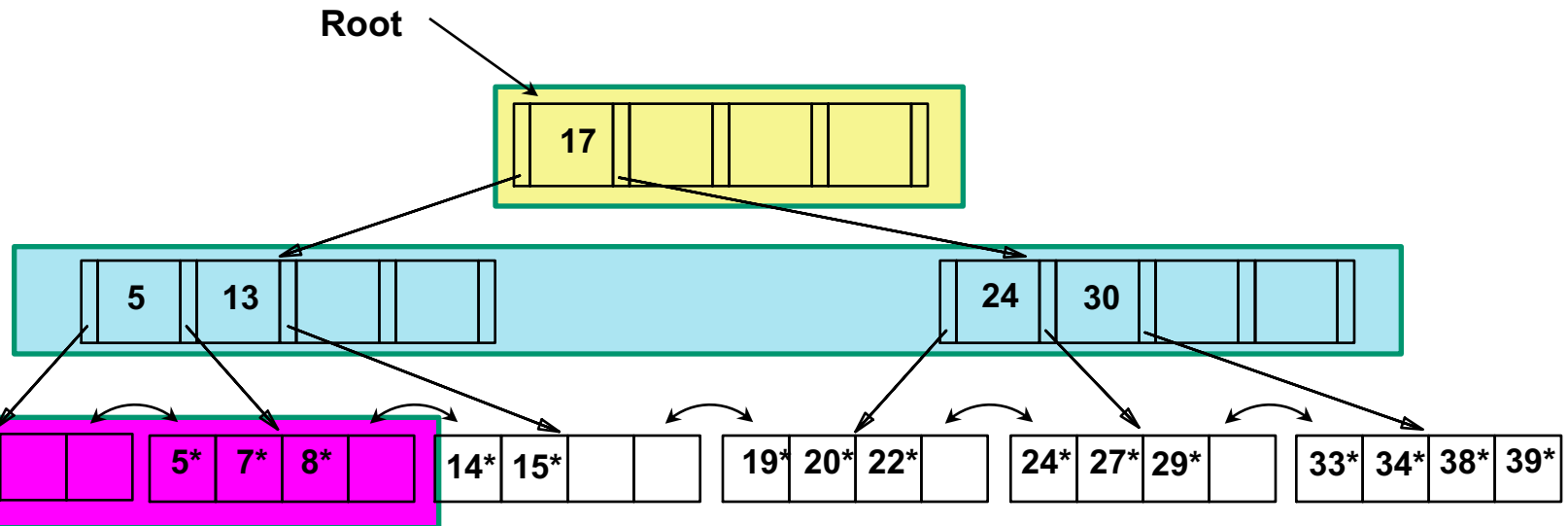
- Understand difference between copy-up and push-up

- Observe how min number of entries (d) is guaranteed in both leaf and index node splits.

**Push up** the index entry to be inserted in parent node



Note that 17 is pushed up and only **appears once** in the index. (Contrast this with a leaf split.)

# Example B+ Tree After Inserting 8*



Notice that root was split, leading to increase in height.

B+ trees grow **bottom-up** dynamically!

# Inserting a Data Entry into a B+ Tree: Summary

- Find correct leaf *L*.

- Put data entry onto *L*.
  - If *L* has enough space, *done*!
  - Else, must *split*  *L (into L and a new node L2)*
    - Redistribute entries evenly, put middle key in L2
    - **copy up** middle key.
    - Insert index entry pointing to *L2* into parent of *L*.

- This can happen recursively
  - To split index node, redistribute entries evenly, but **push up** middle key.  (Contrast with leaf splits.)

- Splits "grow" tree; root split increases height.
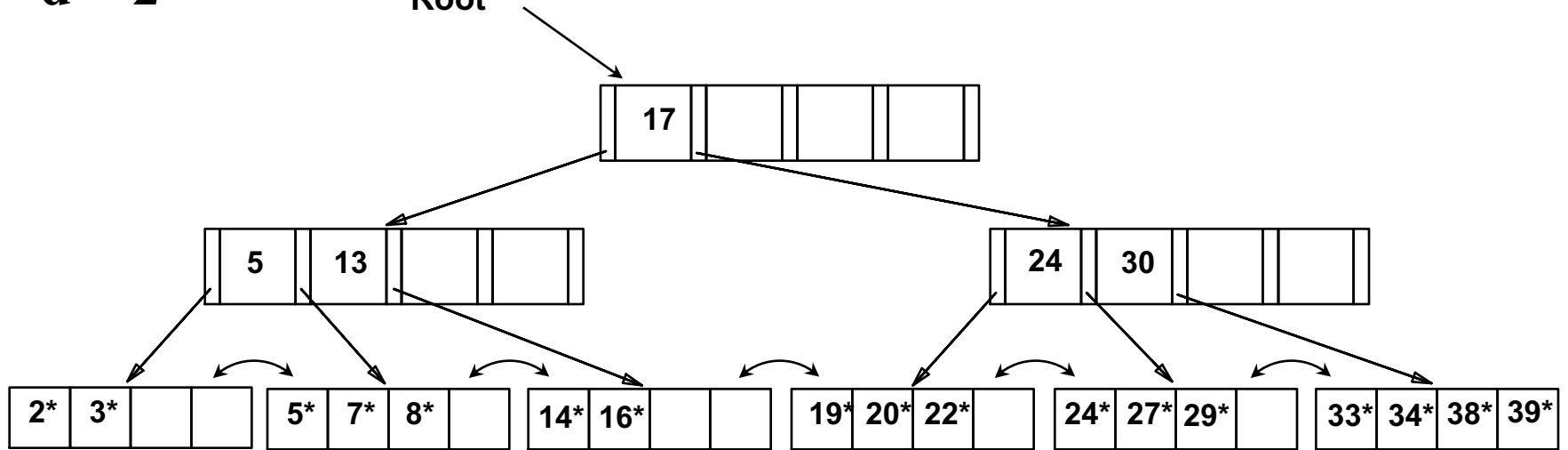  - Tree growth: gets *wider* or *one level taller at top.*

# Deleting a Data Entry from a B+ Tree

- Examine examples first …

# Delete 19* and 20*

**d = 2**
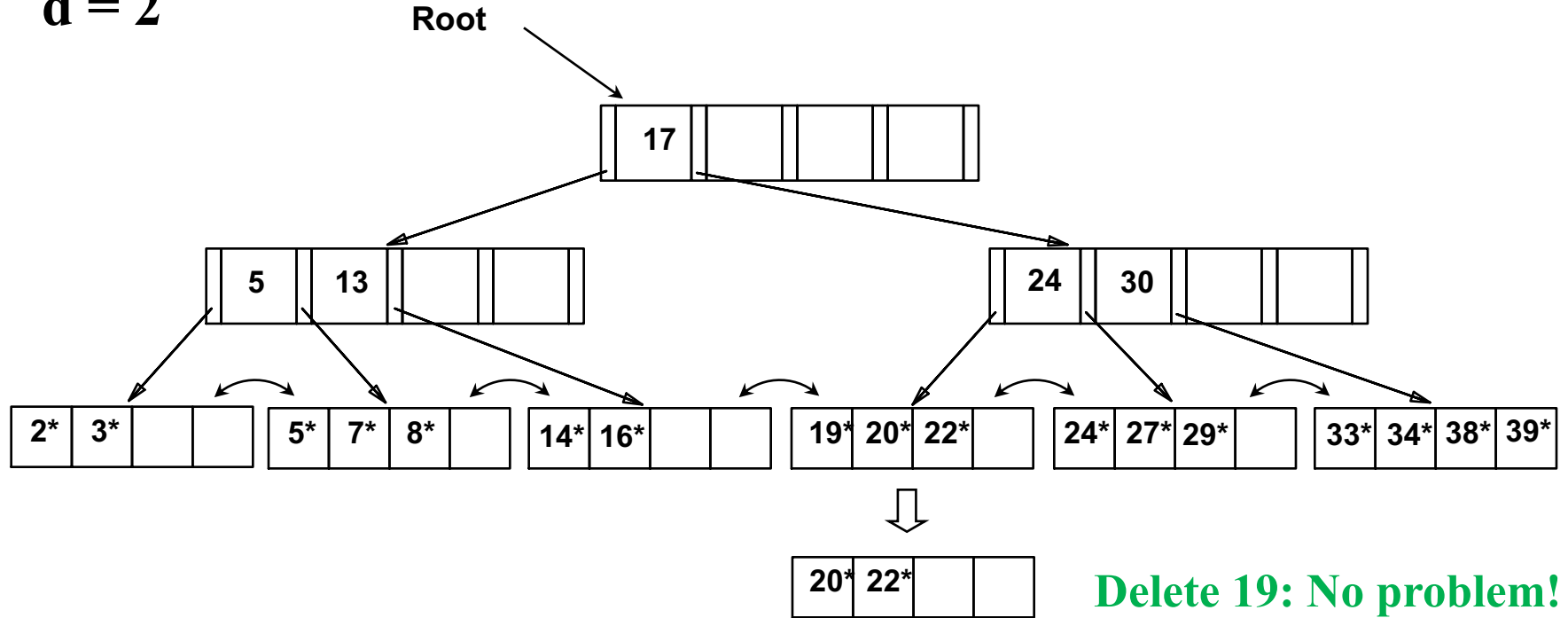
**Root**

| | 17 | | | | | | |

| | 5 | | 13 | | | | |

| | 24 | | 30 | | | | |

| 2* | 3* | | | 5* | 7* | 8* | | 14* | 16* | | | 19* | 20* | 22* | | 24* | 27* | 29* | | 33* | 34* | 38* | 39* |

Leaf nodes:
d ≤ n ≤ 2d
2 ≤ n ≤ 4

# Delete 19* and 20*

**d = 2**

Root

| | 17 | | | | | | |
|---|---|---|---|---|---|---|---|

| | 5 | | 13 | | | | |
|---|---|---|---|---|---|---|---|

| | 24 | | 30 | | | | |
|---|---|---|---|---|---|---|---|

| 2* | 3* | | |
|---|---|---|---|

| 5* | 7* | 8* | |
|---|---|---|---|

| 14* | 16* | | |
|---|---|---|---|

| 19* | 20* | 22* | |
|---|---|---|---|

| 24* | 27* | 29* | |
|---|---|---|---|

| 33* | 34* | 38* | 39* |
|---|---|---|---|

| 20* | 22* | | |
|---|---|---|---|

**Delete 19: No problem!**

Leaf nodes:

$d \leq n \leq 2d$

$2 \leq n \leq 4$

# Delete 19* and 20*

d = 2

Root

17

5 13

24 30

2* 3* | 5* 7* 8* | 14* 16* | 19* 20* 22* | 24* 27* 29* | 33* 34* 38* 39*

20* 22*

**Delete 19: No problem!**

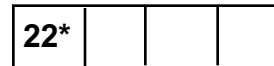Leaf nodes:
d ≤ n ≤ 2d
2 ≤ n ≤ 4

22*

**Delete 20:**
**Leaf node underflows!**

# Delete 19* and 20*

d = 2



**Leaf node underflows!**

**When a leaf node underflows:**

Two options (try in order):

1- Redistribute among sibling nodes evenly, and if this is not possible,

2- Merge nodes

# Delete 19* and 20*

**d = 2**

Root

17

5 | 13

24 | 30

2* | 3*

5* | 7* | 8*

14* | 16*

22*

24* | 27* | 29*

33* | 34* | 38* | 39*

## Redistribute for leaf nodes:

A **sibling** of a node is the node that is <u>adjacent</u> (immediate to left or right) to it, and has the <u>same parent</u>

32

# Delete 19* and 20*

**d = 2**

Root

| 17 | | | |

| 5 | 13 | | |

| 24 | 30 | | |

| 2* | 3* | | |

| 5* | 7* | 8* | |

| 14* | 16* | | |

| 22* | | | |

| 24* | 27* | 29* | |

| 33* | 34* | 38* | 39* |

↑
**This node is Not a sibling!**
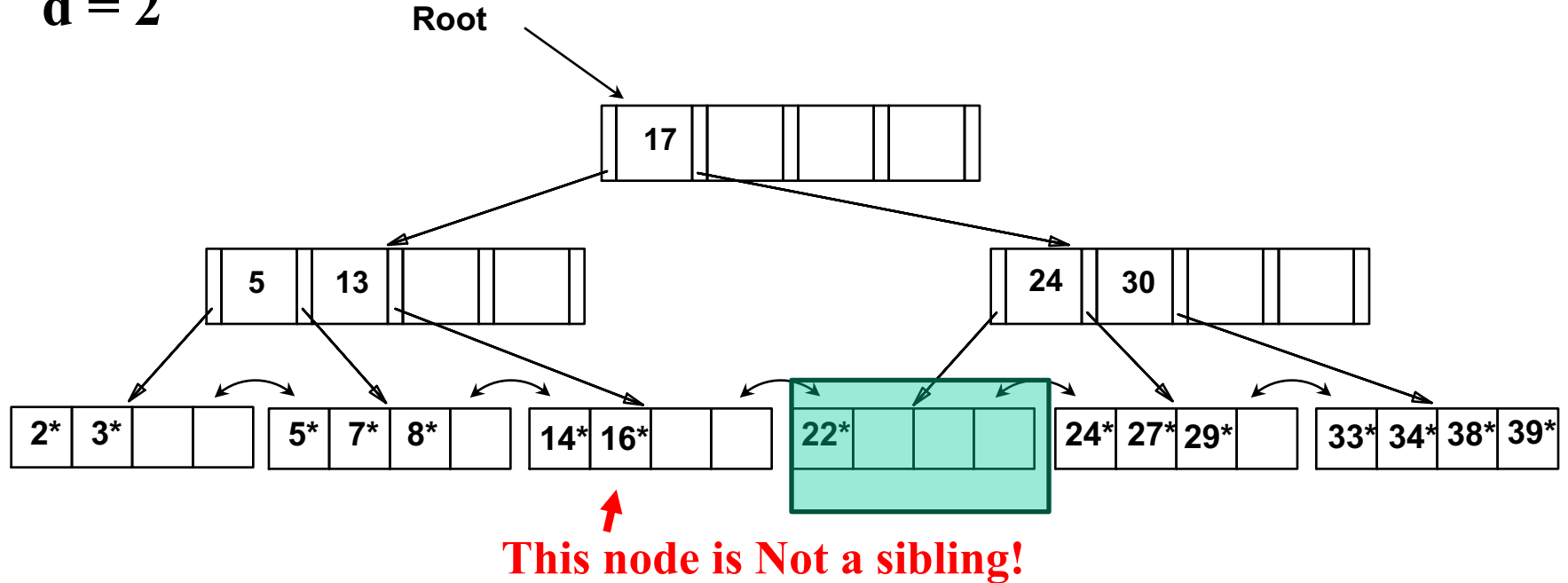
**Redistribute for leaf nodes:**

A **sibling** of a node is the node that is <u>adjacent</u> (immediate to left or right) to it, and has the <u>same parent</u>

33

# Delete 19* and 20*

d = 2

Root

```
                              ┌──┬──┬──┬──┐
                              │17│  │  │  │
                              └──┴──┴──┴──┘
             ┌──┬──┬──┬──┐              ┌──┬──┬──┬──┐
             │ 5│13│  │  │              │24│30│  │  │
             └──┴──┴──┴──┘              └──┴──┴──┴──┘
```

```
┌──┬──┬──┬──┐  ┌──┬──┬──┬──┐  ┌──┬──┬──┬──┐  ┌──┬──┬──┬──┐  ┌──┬──┬──┬──┐  ┌──┬──┬──┬──┐
│2*│3*│  │  │  │5*│7*│8*│  │  │14*│16*│ │ │  │22*│  │  │  │  │24*│27*│29*│ │  │33*│34*│38*│39*│
└──┴──┴──┴──┘  └──┴──┴──┴──┘  └──┴──┴──┴──┘  └──┴──┴──┴──┘  └──┴──┴──┴──┘  └──┴──┴──┴──┘
```

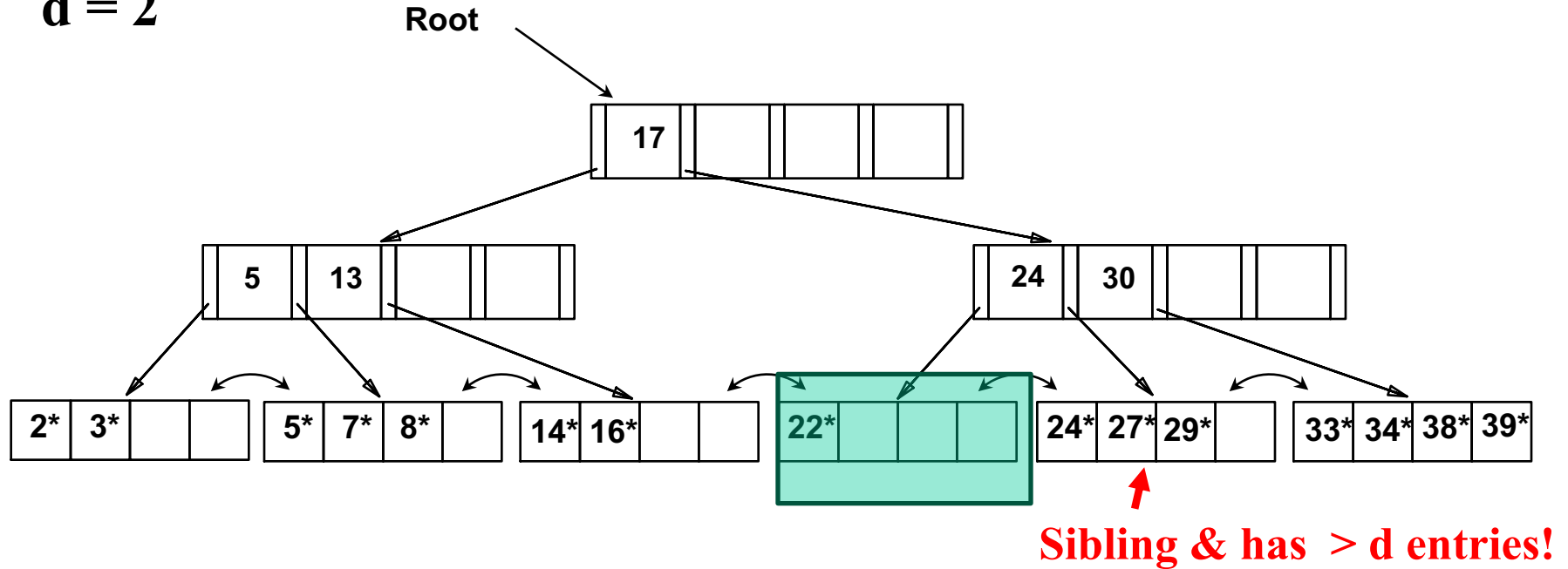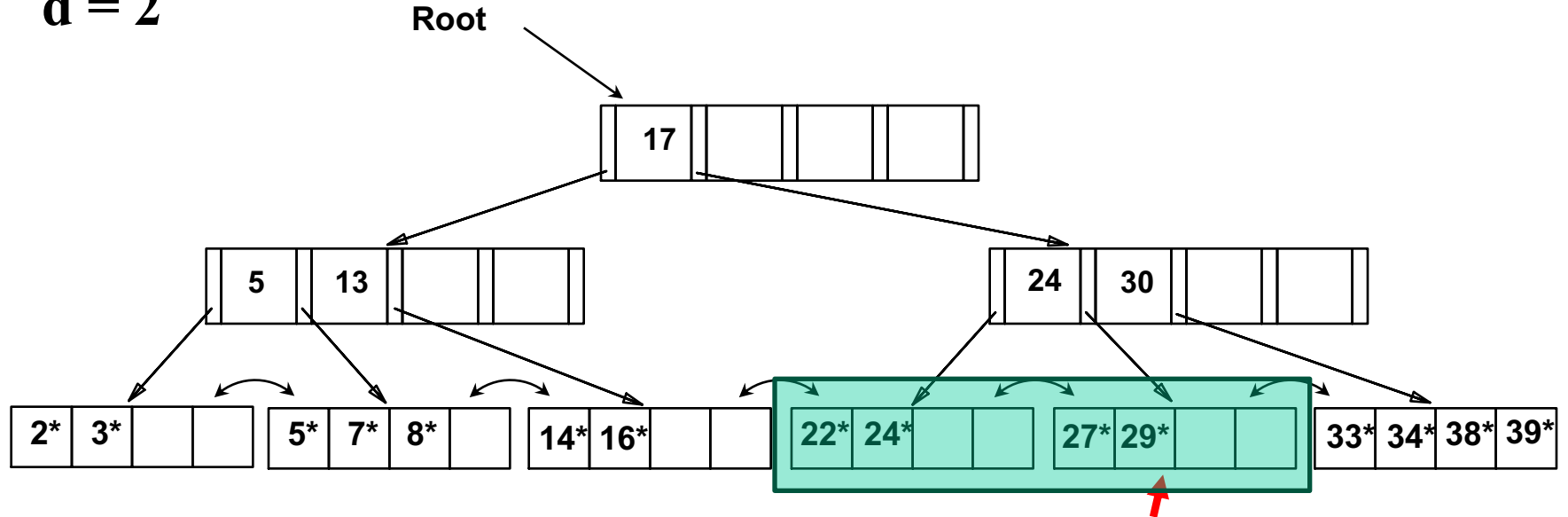**Sibling & has > d entries!**

# Redistribute for leaf nodes:

A **sibling** of a node is the node that is <u>adjacent</u> (immediate to left or right) to it, and has the <u>same parent</u>

34

# Delete 19* and 20*

**d = 2**

Root

| 17 | | | |
|----|--|--|--|

| 5 | 13 | | |
|---|----|--|--|

| 24 | 30 | | |
|----|----|--|--|

| 2* | 3* | | |
|----|----|--|--|

| 5* | 7* | 8* |
|----|----|----|

| 14* | 16* | |
|-----|-----|--|

| 22* | 24* | | |
|-----|-----|--|--|

| 27* | 29* | | |
|-----|-----|--|--|

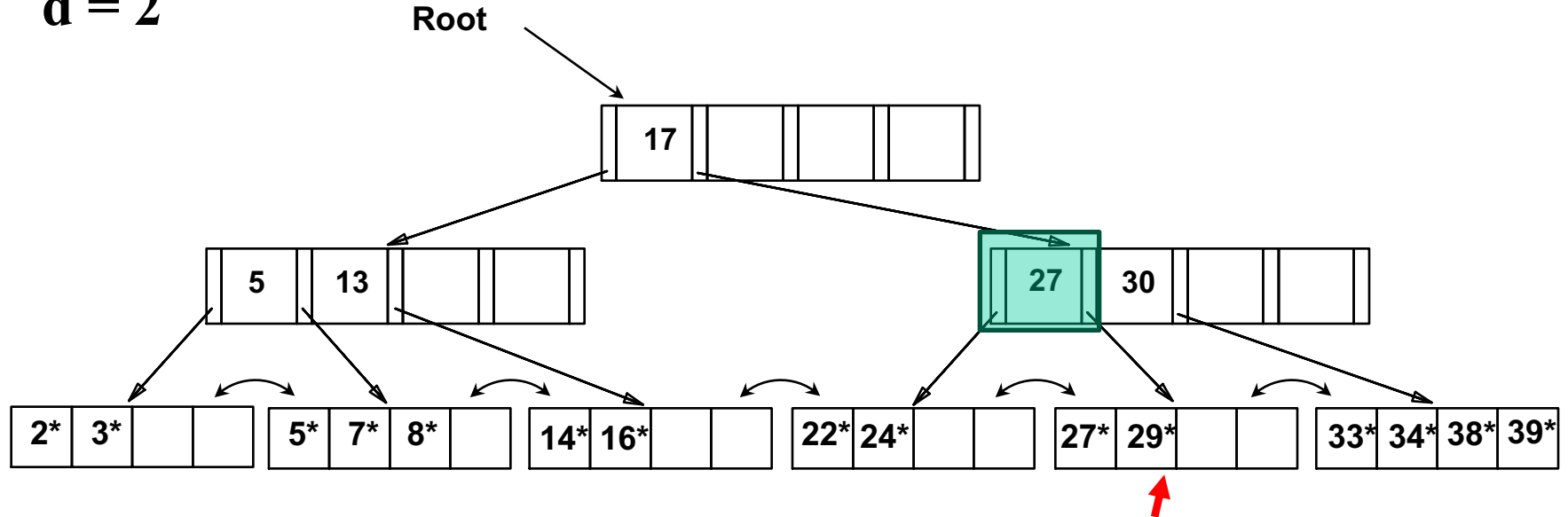| 33* | 34* | 38* | 39* |
|-----|-----|-----|-----|

**Sibling & has > d entries!**

**Redistribute for leaf nodes:**

A **sibling** of a node is the node that is <u>adjacent</u> (immediate to left or right) to it, and has the <u>same parent</u>

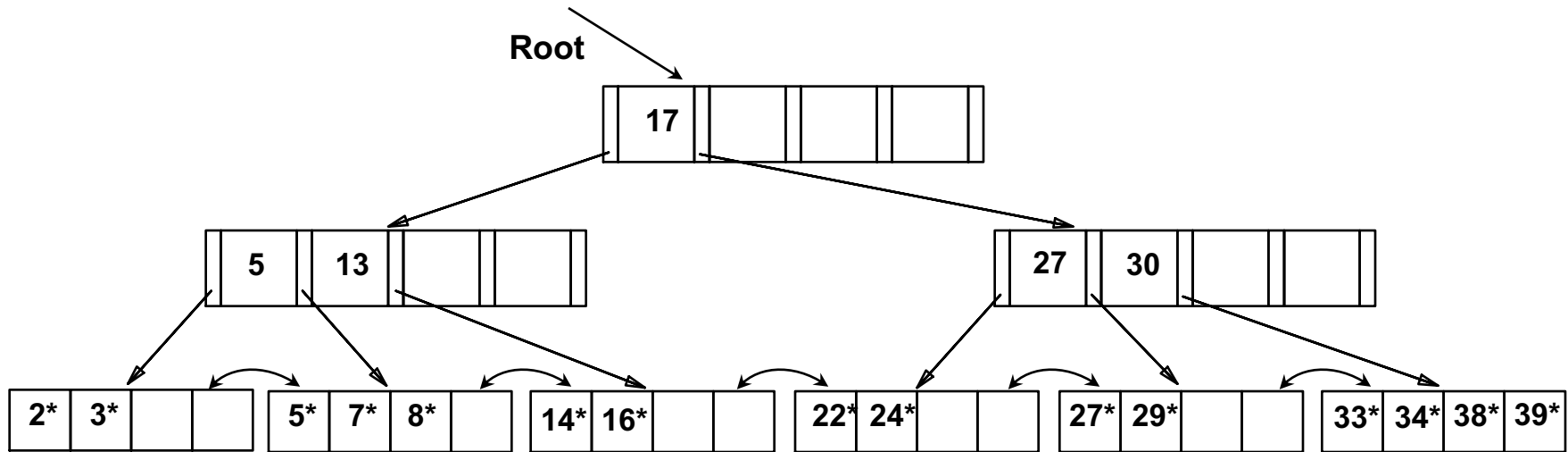1) Redistribute among siblings

# Delete 19* and 20*

**d = 2**



**Sibling & has > d entries!**

**Redistribute for leaf nodes:**

A **sibling** of a node is the node that is <u>adjacent</u> (immediate to left or right) to it, and has the <u>same parent</u>

1)  Redistribute among siblings
2)  COPY-UP (Update) the middle key as the search key

36

# Deleting 19* and 20* (cont.)
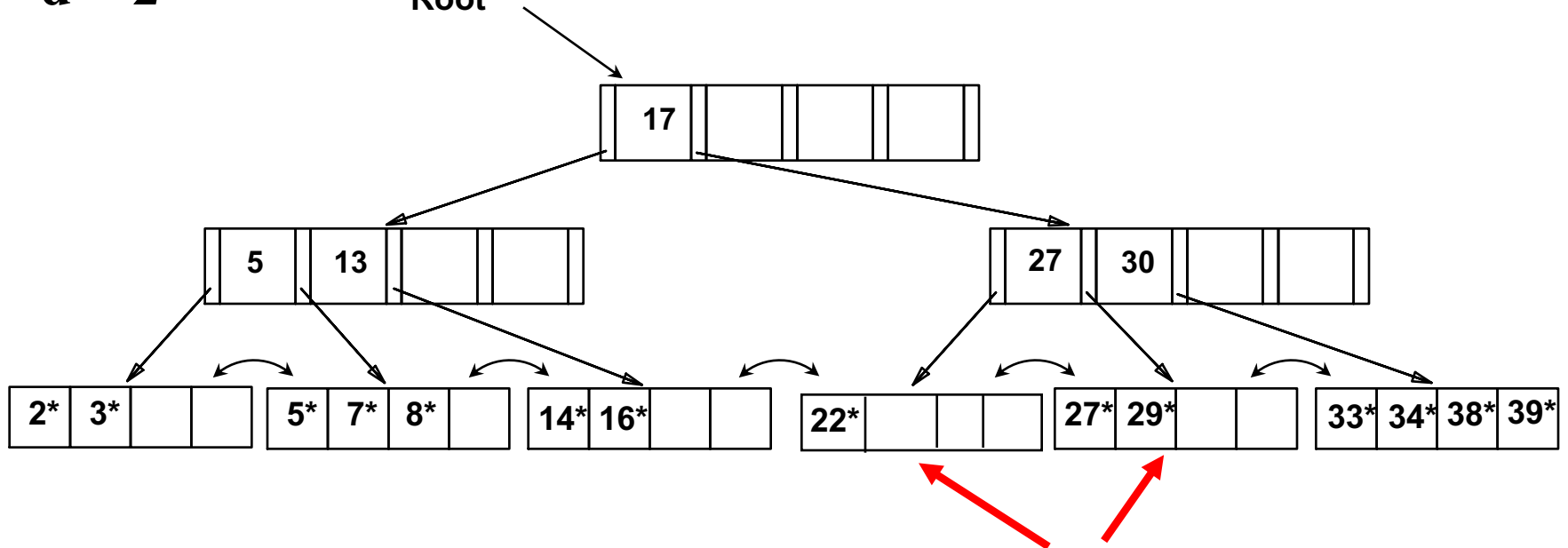
**Root**

```
                    17
        5   13              27   30

2* 3*    5* 7* 8*   14* 16*    22* 24*    27* 29*    33* 34* 38* 39*
```

- Notice how 27 is *copied up*.
- But can we move it up?
- Now we want to delete 24
- Underflow again!

# Delete 24*

d = 2

Root

```
                                    [ 17 |   |   |   ]

        [ 5 | 13 |   |   ]                          [ 27 | 30 |   |   ]

[2*|3*|  |  ]  [5*|7*|8*|  ]  [14*|16*|  |  ]  [22*|  |  |  ]  [27*|29*|  |  ]  [33*|34*|38*|39*]
```

**Option (1) Not Applicable here!**
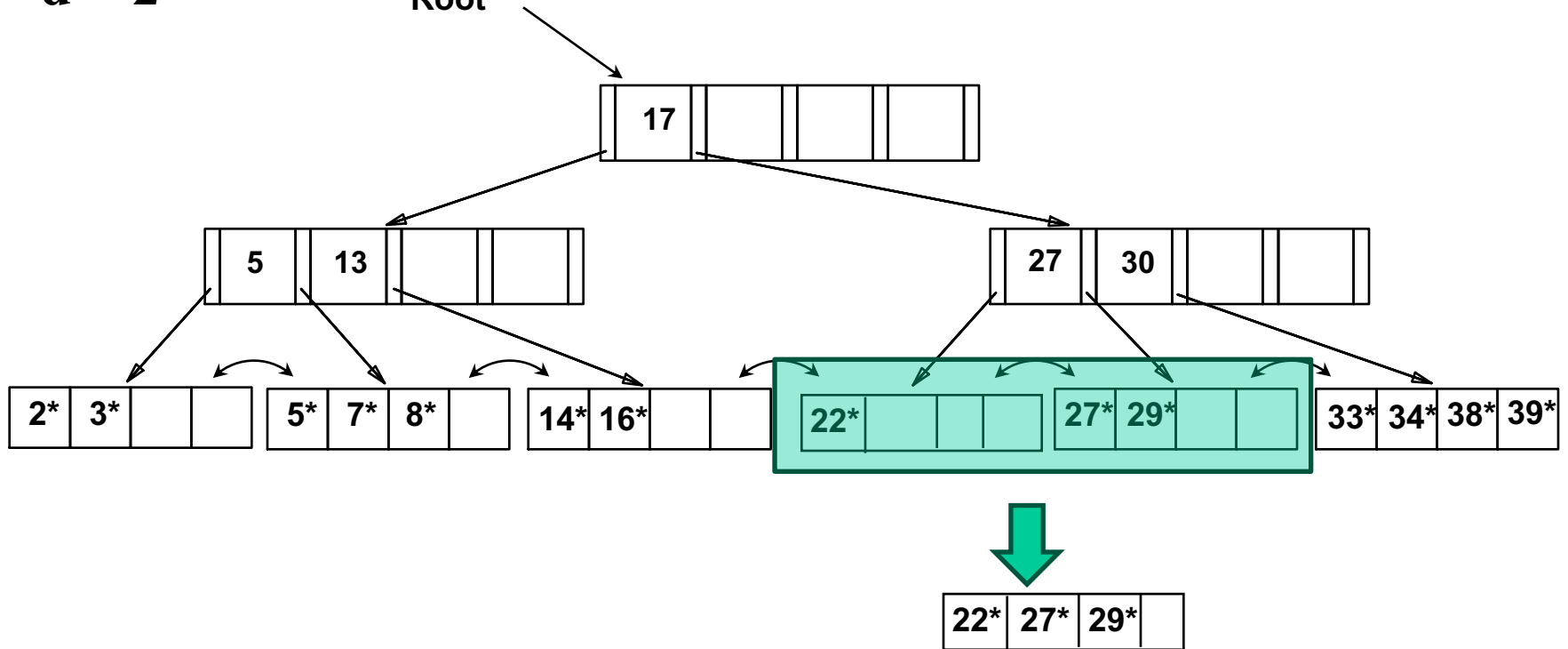
**When a <u>leaf node</u> underflows:**

Two options (try in order):

1- Redistribute among sibling nodes evenly, and if this is not possible,
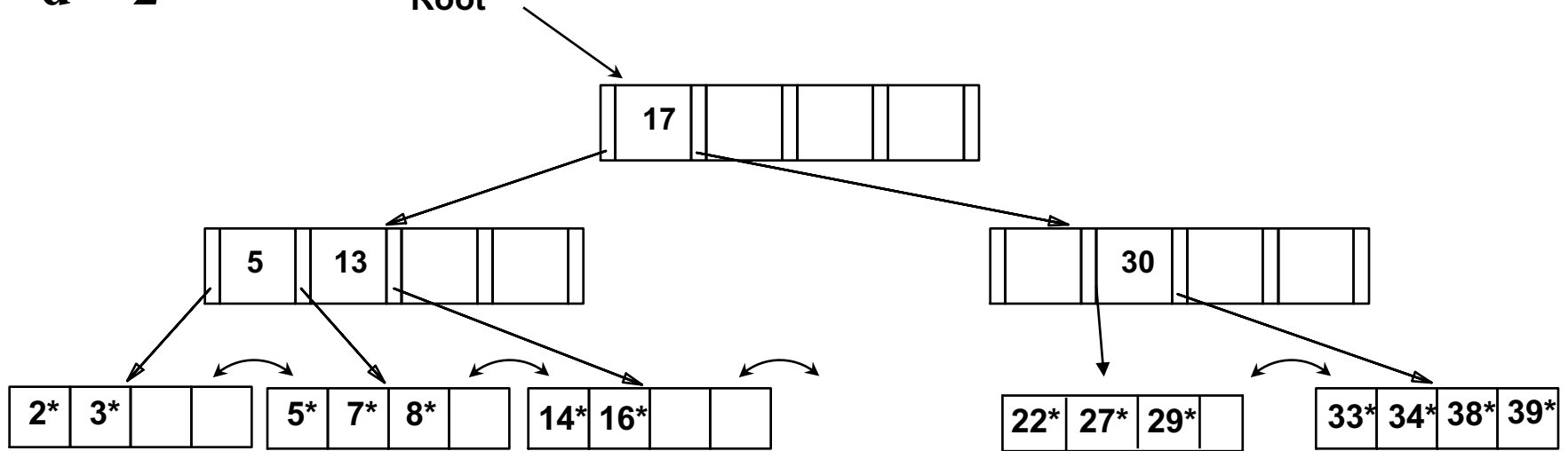
2- Merge nodes

38

# Delete 24*

d = 2



**Option (2): Merge leaf nodes**

Step 1: Merge leaf nodes

# Delete 24*

**d = 2**



**Option (2): Merge leaf nodes**

Step 1: Merge leaf nodes
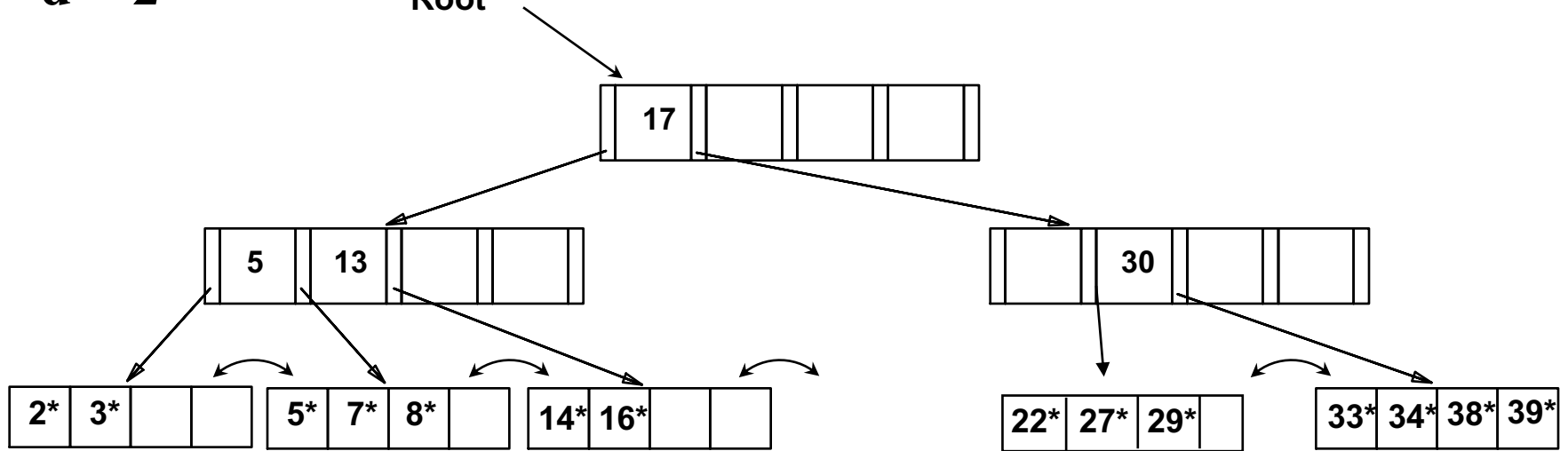Step 2: Remove the search key entry and pointer to the discarded node

# Delete 24*

**d = 2**

Root



| 17 | | | |

| 5 | 13 | | |

| 30 | | |

| 2* | 3* | | |

| 5* | 7* | 8* | |

| 14* | 16* | | |

| 22* | 27* | 29* | |

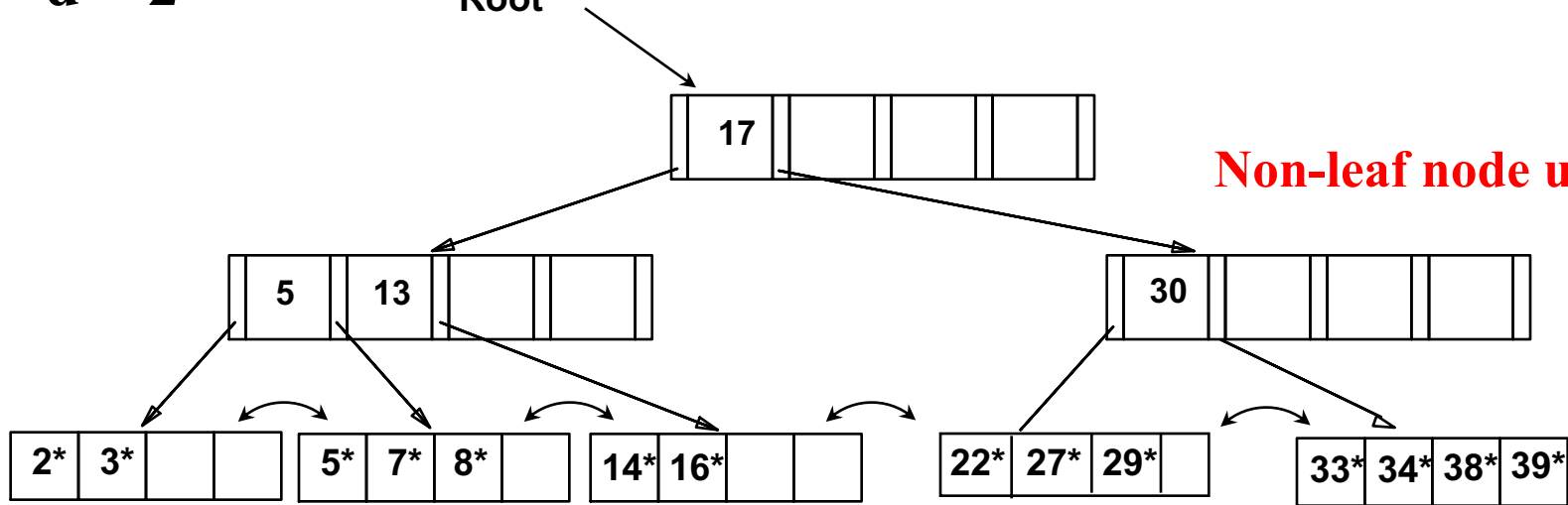| 33* | 34* | 38* | 39* |

**Is it good like this?**

**Option (2): Merge leaf nodes**

Step 1: Merge leaf nodes
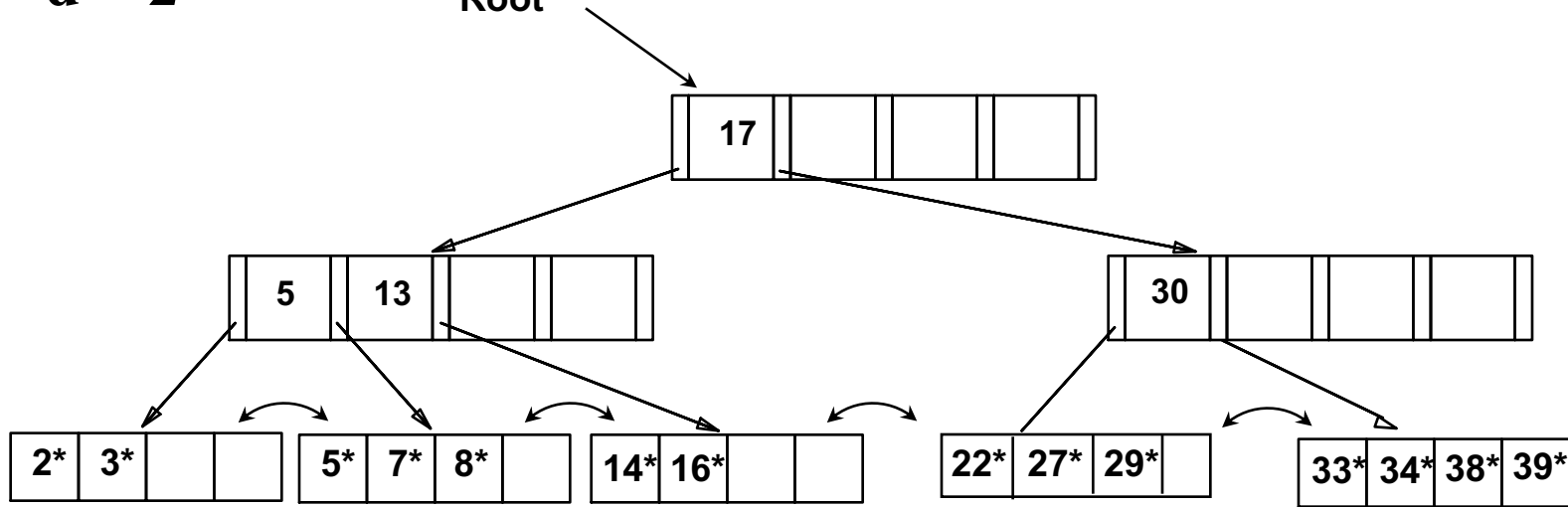Step 2: Remove the search key entry and pointer to the discarded node

# Delete 24*

**d = 2**

Root

17

**Non-leaf node underflows!**

5   13

30

| 2* | 3* | | |

| 5* | 7* | 8* | |

| 14* | 16* | | |

| 22* | 27* | 29* | |

| 33* | 34* | 38* | 39* |

⇩

# Delete 24*

**d = 2**

Root



| | 17 | | | | | |

| | 5 | 13 | | | |     | | 30 | | | | |

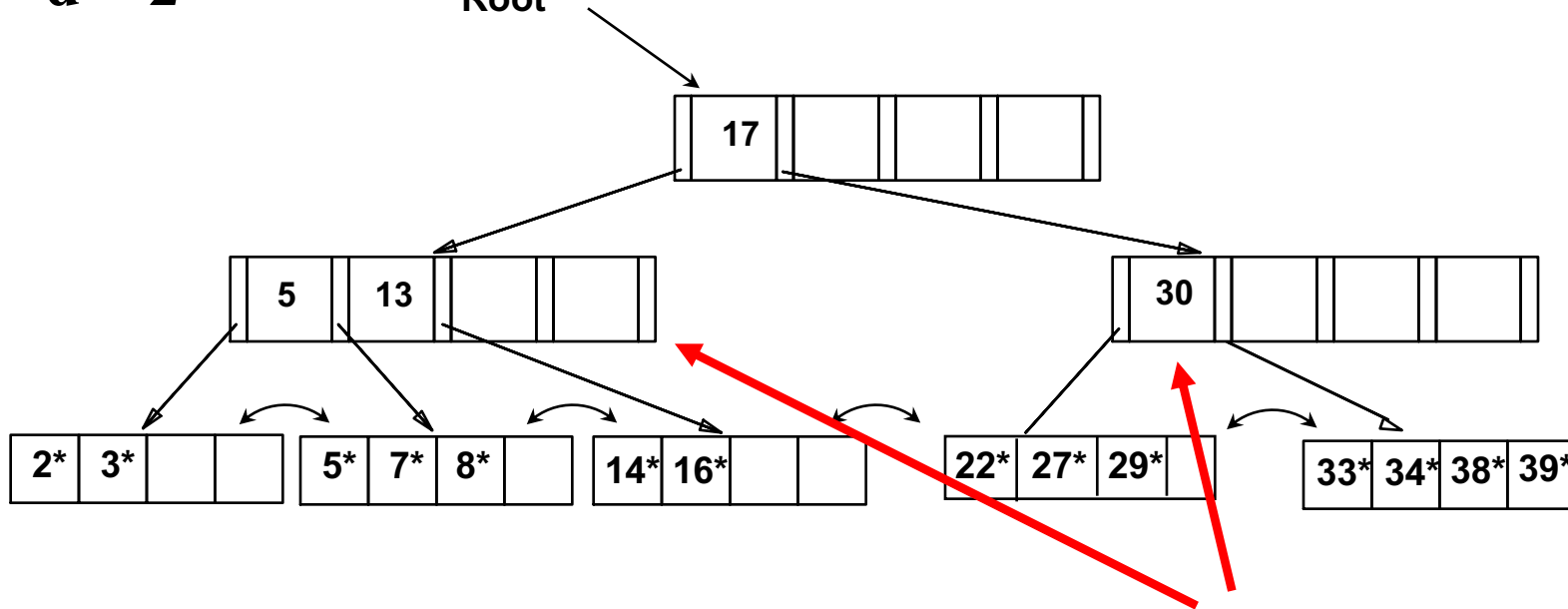| 2* | 3* | | |   | 5* | 7* | 8* | |   | 14* | 16* | | |     | 22* | 27* | 29* | |   | 33* | 34* | 38* | 39* |

**When a <u>non-leaf node</u> underflows:**

Two options (try in order):

1- Redistribute among sibling nodes evenly evenly, and if this is not possible,

2- Merge nodes

# Delete 24*
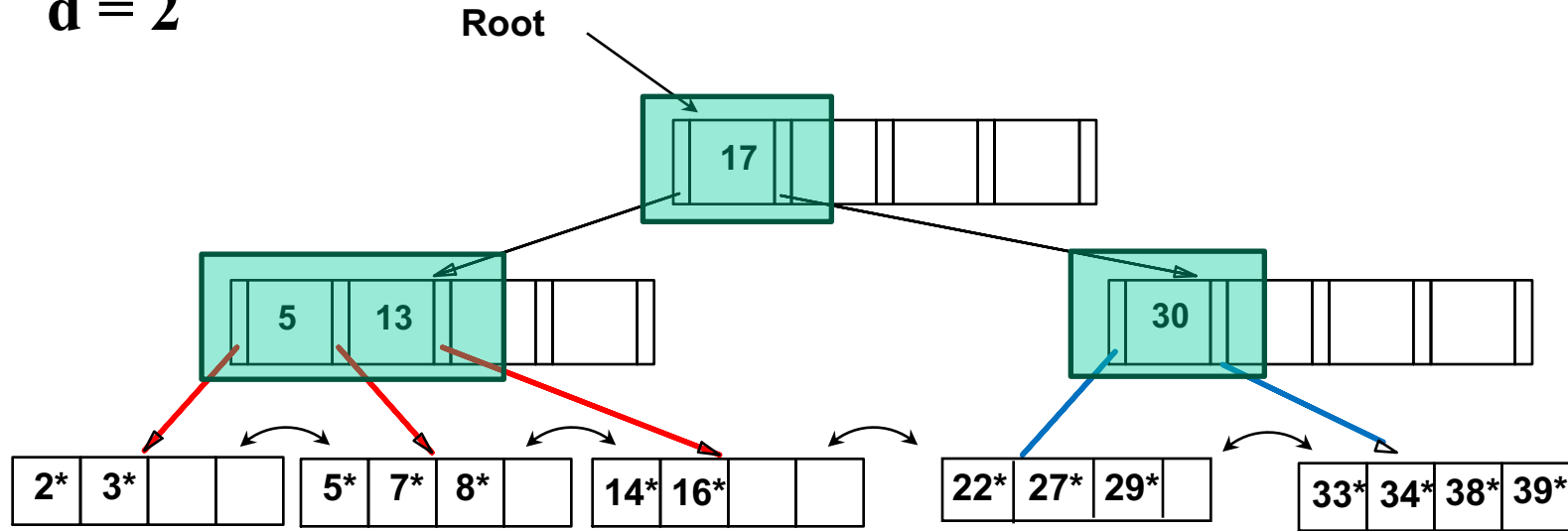
**d = 2**



**Option (1) Not Applicable here!**

**When a <u>non-leaf node</u> underflows:**

Two options (try in order):

1- Redistribute among sibling nodes evenly evenly, and if this is not possible,

2- Merge nodes

# Delete 24*

d = 2

Root



17

5    13

30

2*  3*        5*  7*  8*        14* 16*        22* 27* 29*        33* 34* 38* 39*

**Option (2): Merge**
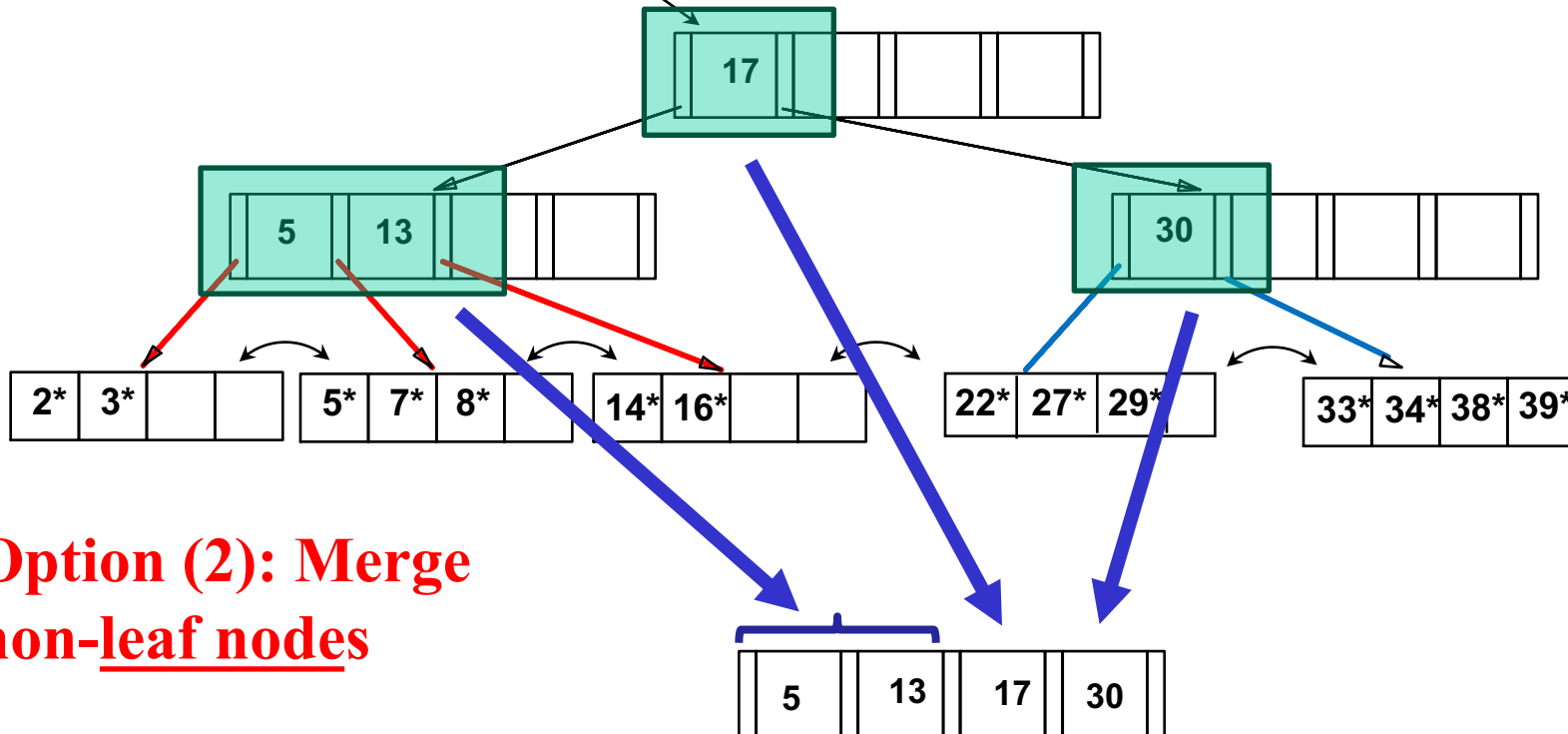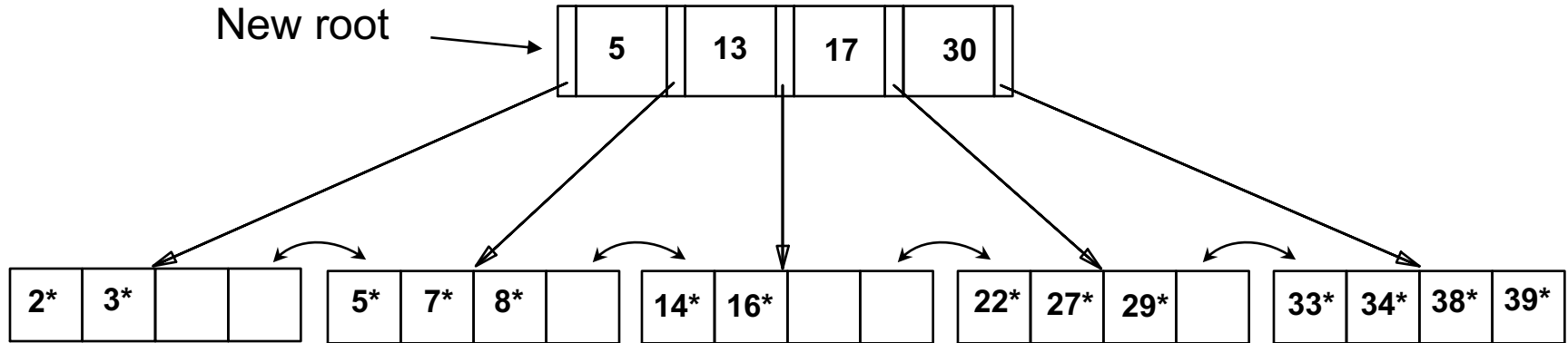**non-leaf nodes**

Merge:
- Entries in first non-leaf node (together with pointers),
- **PULL DOWN** the splitting search key,
- followed by the entries in the second non-leaf node
  (together with pointers)

# Delete 24*

**d = 2**

Root



**17**

**5**  **13**

**30**

| 2* | 3* | | |
| 5* | 7* | 8* | |
| 14* | 16* | | |
| 22* | 27* | 29* |
| 33* | 34* | 38* | 39* |

**Option (2): Merge non-leaf nodes**

| | 5 | | 13 | | 17 | | 30 | |

Merge:
- Entries in first non-leaf node (together with pointers),
- **PULL DOWN** the splitting search key,
- followed by the entries in the second non-leaf node (together with pointers)

# Delete 24*

New root → | | 5 | 13 | 17 | 30 | |

| 2* | 3* | | | | 5* | 7* | 8* | | | 14* | 16* | | | | 22* | 27* | 29* | | | 33* | 34* | 38* | 39* |

# Deleting a Data Entry from a B+ Tree: Summary

- Start at root, find leaf *L* where entry belongs.
- Remove the entry.
  - If L is at least half-full, *done!*
  - If L has only **d-1** entries,
    - Try to re-distribute, borrowing from *sibling (adjacent node with same parent as L)*.
    - If re-distribution fails, *merge* L and sibling.
- If merge occurred, must delete entry (pointing to *L* or sibling) from parent of *L*.
- Merge could propagate to root, decreasing height.

# Non-leaf Node Redistribution

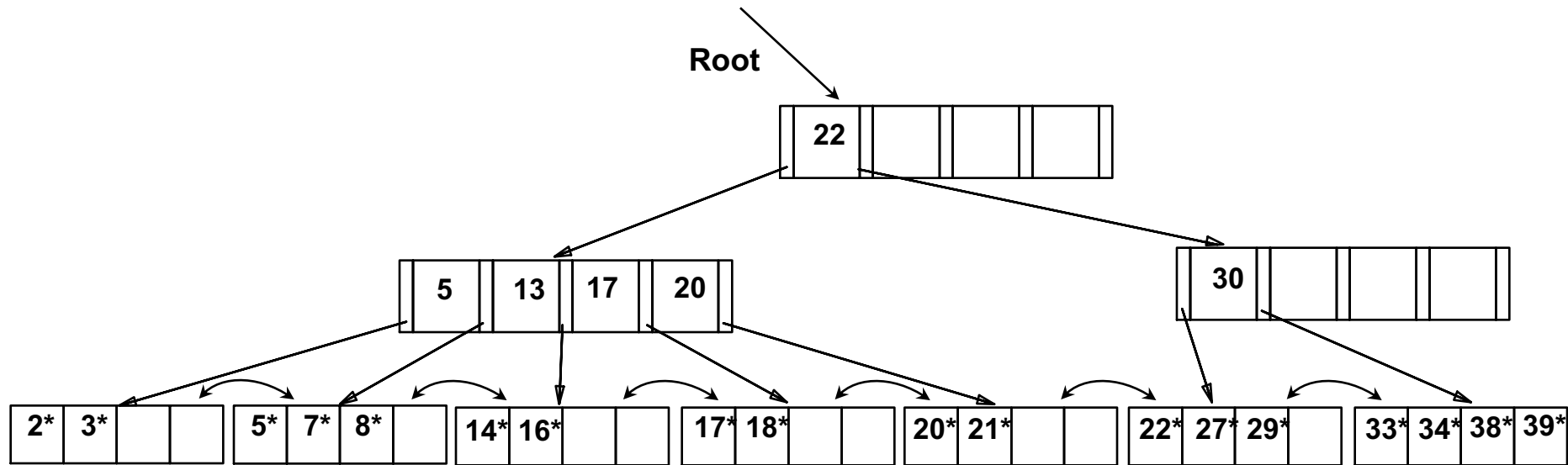**When a <u>non-leaf node</u> underflows:**

Two options (try in order):

1- Redistribute among nodes evenly, and if this is not possible,

2- Merge nodes

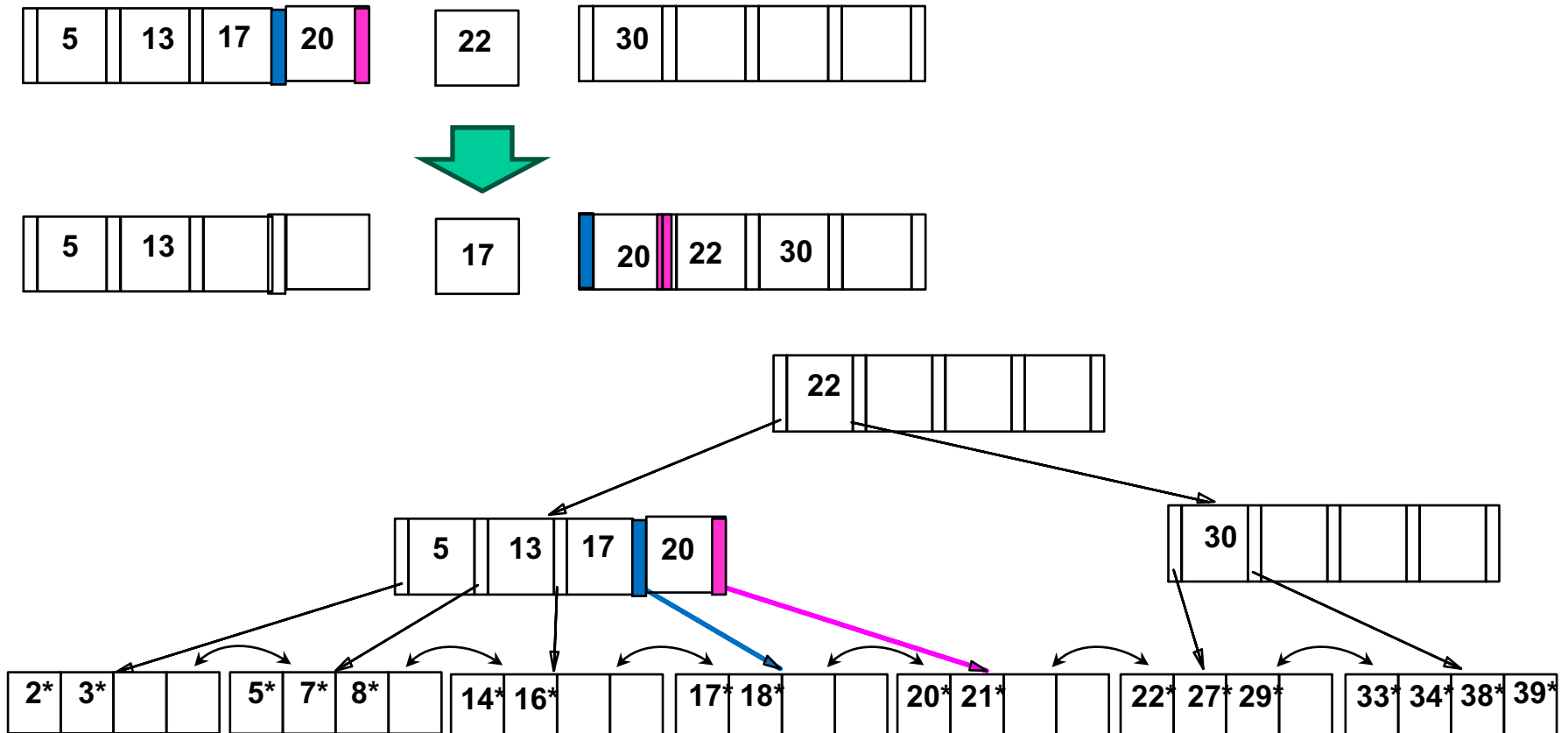We have already seen an example for the second case.

How about the first case!

# Example of Non-leaf Re-distribution

- Tree is shown below *during deletion* of 24*. (What could be a possible initial tree?)
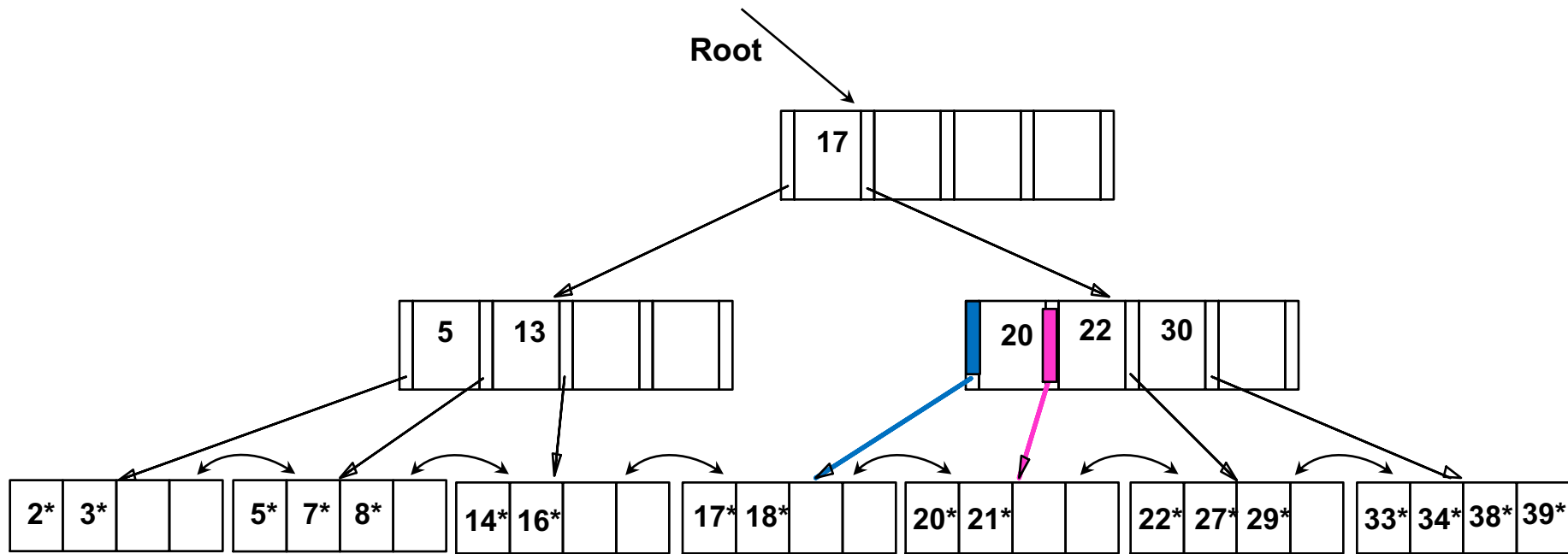- In contrast to previous example, can re-distribute entry from left child of root to right child.

# Re-distribution

- Intuitively, entries are re-distributed by `*pushing through*' the splitting entry in the parent node.

- Consider all search keys together

| 5 | 13 | 17 | | 20 | |

| 22 |

| 30 | | | |

| 5 | 13 | | |

| 17 |

| | 20 | 22 | 30 | |

| 22 | | | |

| 5 | 13 | 17 | 20 | |

| 30 | | | |

| 2* | 3* | | | 5* | 7* | 8* | | 14* | 16* | | | 17* | 18* | | | 20* | 21* | | | 22* | 27* | 29* | | 33* | 34* | 38* | 39* |

# After Re-distribution

**Root**

| 17 | | | |

| 5 | 13 | | |

| | 20 | 22 | 30 | |

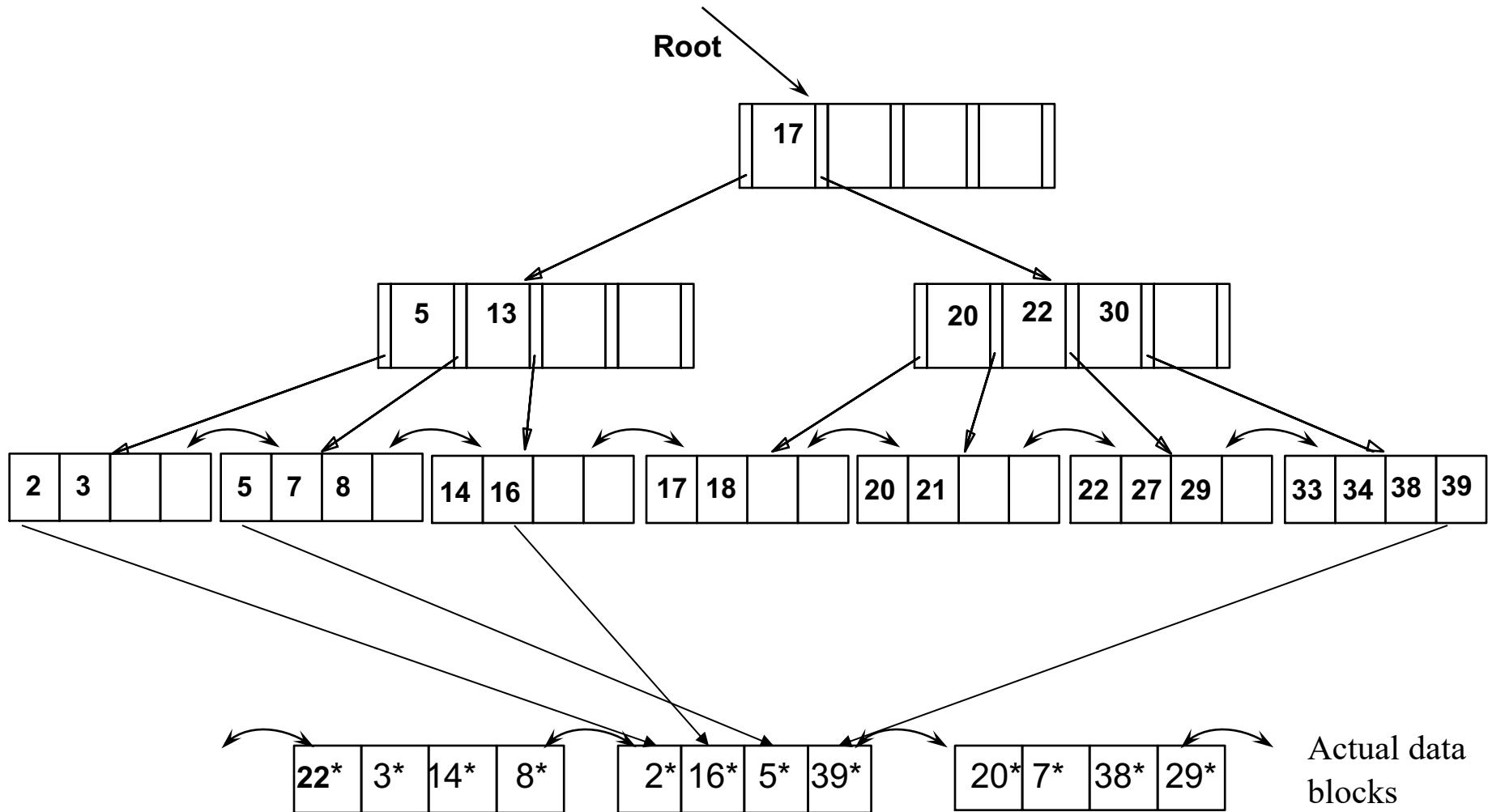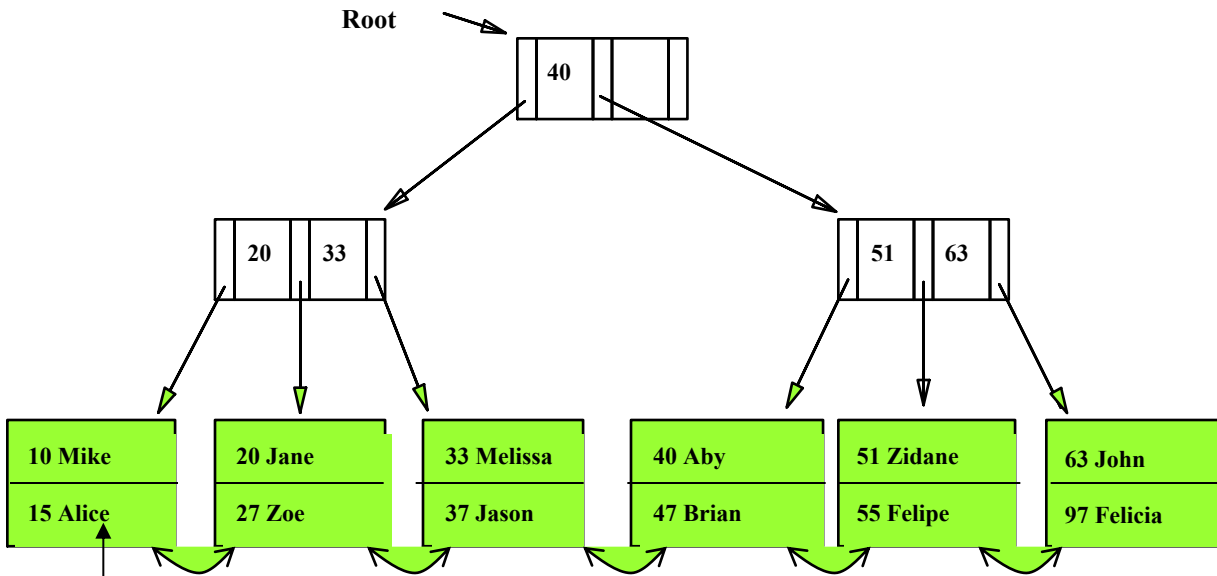| 2* | 3* | | |   | 5* | 7* | 8* | |   | 14* | 16* | | |   | 17* | 18* | | |   | 20* | 21* | | |   | 22* | 27* | 29* | |   | 33* | 34* | 38* | 39* |

# Primary vs Secondary Index

- Note: We were assuming the data items were in sorted order
  - This is called *primary/clustered B+tree* index
- *Secondary B+tree* index:
  - Built on an attribute that the file is not sorted on.
- Can have many different indexes on the same file.

# A Secondary B+-Tree index

**Root**

| 17 | | | |

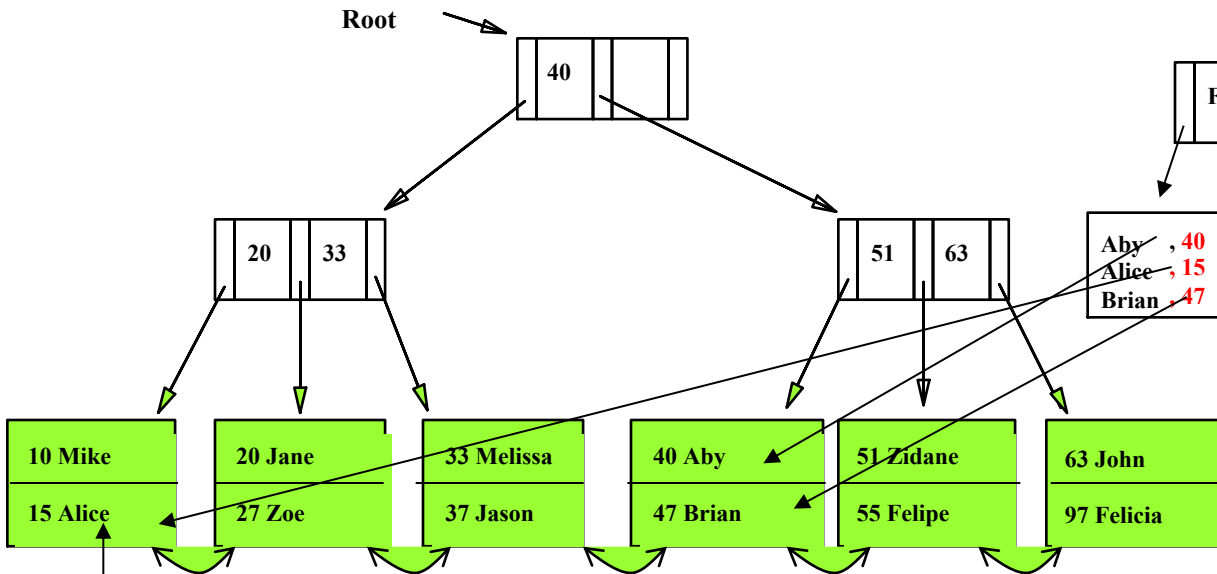| 5 | 13 | | |

| 20 | 22 | 30 | |

| 2 | 3 | | |

| 5 | 7 | 8 | |

| 14 | 16 | | |

| 17 | 18 | | |

| 20 | 21 | | |

| 22 | 27 | 29 | |

| 33 | 34 | 38 | 39 |

| 22* | 3* | 14* | 8* |

| 2* | 16* | 5* | 39* |

| 20* | 7* | 38* | 29* |

Actual data blocks

# A file organized as (or, has) a **Primary B+-Tree** index on *ssn*

**Root**

| | 40 | | |

| | 20 | 33 | |

| | 51 | 63 | |

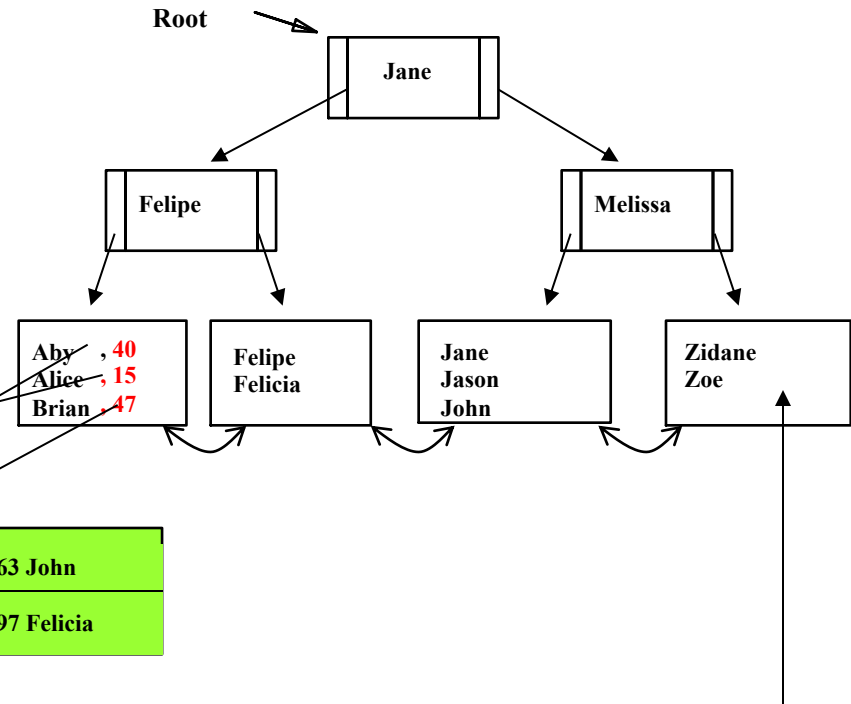| 10 Mike | | 20 Jane | | 33 Melissa | | 40 Aby | | 51 Zidane | | 63 John |
| 15 Alice | | 27 Zoe | | 37 Jason | | 47 Brian | | 55 Felipe | | 97 Felicia |

• As **15***, we store the **actual data record** with key value 15 (**Alternative-1**)
• In this case, the leaf nodes can be larger, say a few blocks (pages)

# A file organized as (or, has) a **Primary B+-Tree** index on *ssn*

## The same file also has a **Secondary B+-Tree** index on *name*

**Root**

| 40 | |

| 20 | 33 | | 51 | 63 |

| 10 Mike | 20 Jane | 33 Melissa | 40 Aby | 51 Zidane | 63 John |
| 15 Alice | 27 Zoe | 37 Jason | 47 Brian | 55 Felipe | 97 Felicia |

**Root**

| Jane | |

| Felipe | | Melissa | |

| Aby , **40** |  | Felipe | | Jane | | Zidane |
| Alice , **15** |  | Felicia | | Jason | | Zoe |
| Brian , **47** |  | | | John | | |

- As **15\***, we store the **actual data record** with key value 15 (**Alternative-1**)
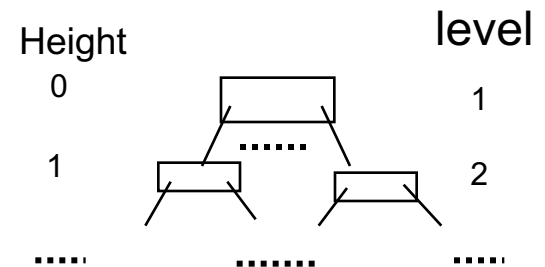- In this case, the leaf nodes can be larger, say a few blocks (pages)

- We have k* as <key, rid> (**Alternative 2**)
- We can have rid's as pointers, or use PK as rid. We show both above

# Cost for searching a value in B+ tree

- Assumptions:
  - Each interior node is a disk block
  - Each leaf node is also a disk block and data entries (K*) are of the form <key, ptr>. There are D data entries.
  - Let F be the average number of pointers in a node (for internal nodes, it is called *fanout,* **i.e.,** avg. number of children )
- Observe: Let H be the height of the B+ tree: we need to read H+1 nodes (blocks) to reach a data entry in a leaf node
- How do we find H?
  - Level 1 = 1 page = $F^0$ page
  - Level 2 = F pages = $F^1$ pages
  - Level 3 = F * F pages = $F^2$ pages
  - Level H+1 = …….. = $F^H$ pages (i.e., leaf nodes)
  - F pointers → F-1 keys, so there must be D/(F-1) leaf nodes
  - D/(F-1) = $F^H$. That is, $H = \log_F(\frac{D}{F-1})$

Height 0, level 1

Height 1, level 2

# B+ Trees in Practice

- Typical order: 100.  Typical fill-factor: 66%.
  - average fanout = 133 (i.e, # of pointers in internal node)
- Can often hold top levels in buffer pool:
  - Level 1 =            1 page  =     8 Kbytes
  - Level 2 =      133 pages =     1 Mbyte
  - Level 3 = 17,689 pages = 133 MBytes
- Suppose there are 1,000,000,000 data entries.
  - $H = \log_{133}(1000000000/132) < 4$
  - The cost is reading H+1 = 5 pages

# Cost Computation: Another Example

Leaves would store the **actual records**

- A <u>primary B+ tree index</u> on key field giftID.

- 2.500.000 gift records, each record: 400 bytes.

- giftID: 12 bytes, address pointer: 4 bytes

- A <u>bucket</u> can hold 500 records
  - So we have larger leaf nodes (called **buckets**), as we store actual records
  - No claim for interior nodes, assume each is a block!

- B+ tree will have a fill factor of 50% [min occupancy]

- B (block size): 1600

- s: 10 ms, r: 5 ms, btt: 1 ms.

# a) No of index nodes and their total size

We need to find i) fanout of the nodes, and ii) no of leaves.

**i) fanout**: Assume , **n keys (n+1) ptrs** can fit to an index node:

$n \times 12 + (n+1) \times 4 = 1600$ bytes $\rightarrow 16n = 1596 / 16 \rightarrow n = 99$

So at most 99 keys in a node ( 2d= 99, d (tree order) is floor(99/2) )

Tree fill factor 50%; max 99 keys x 50% = 49 keys

fanout: 49 + 1 = 50 ptrs per node

**ii) no of leaves:**

500 rec/leaf * fill factor (50%) = 250 recs/leaf

2.5M records / 250 = 10000 leaf nodes (i.e., buckets)

# a) No of index nodes and their total size

- Tree height = $\log_{50} 10000 = 3$
- So, there are H+1 = 4 levels

Level 4: 10000 leaf nodes (data buckets)

Level 3: ceil (10000 / 50 ptrs) = 200 nodes

Level 2: ceil (200/50) = 4 nodes

Level 1: ceil (4/50) = 1 node (root)

Index nodes: 1 + 4 + 200 = 205

Total Size: 205 x 1600 bytes

## b) Time cost of reading an arbitrary record

- Three has H=3, so 4 levels
- At the first 3 levels, we fetch index nodes:

3 x (s + r+ btt) = 3x (10 +5+1) = 48 ms

- At the fourth level we fetch the leaf node (data bucket)
  - But how many blocks is a data bucket?
  - (500 recs x 400 bytes/rec) / 1600 = 125 blocks
  - So, cost s + r+ 125 x btt = 10+ 5+ 125 x 1 = 140 ms
- Total cost: 48 + 140 = 188 ms

# c) **Cost of reading all records in sorted manner**

- Reach to leftmost leaf node, as before:
- at the first 3 levels, we fetch index nodes:

$3 \times (s + r + btt) = 3 \times (10 + 5 + 1) = 48$ ms

- Read all the leaf nodes (using doubly linked list pointers)
  - $10000 (s + r + 125 \times btt)$
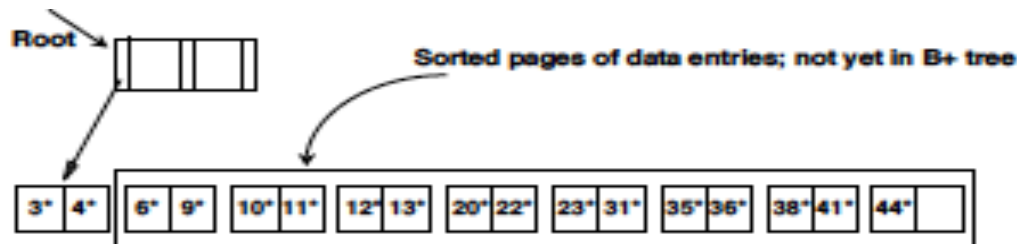- Think: What if this is a secondary B+ tree and we store <key, ptr> pairs at leaf nodes (data buckets)?

# Terminology

- **Blocking Factor:** the number of records which can fit in a leaf node.

- **Fan-out :** the average number of children of an internal node.

- A B+tree index can be used either as a primary index or a secondary index.

  - **Primary index:** determines the way the records are actually stored

  - **Secondary index:** the records in the file are not grouped in blocks according to keys of secondary indexes
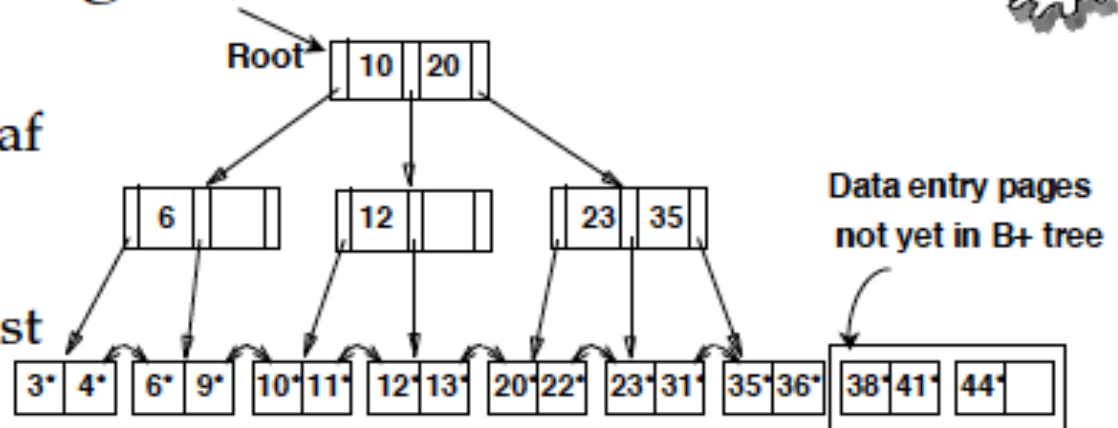
# Bulk Loading of a B+ Tree

- If we have a large collection of records, and we want to create a B+ tree on some field, doing so by repeatedly inserting records is very slow.

- Bulk Loading can be done much more efficiently.

  – Initialization:  Sort all data entries, insert pointer to first (leaf) page in a new (root) page
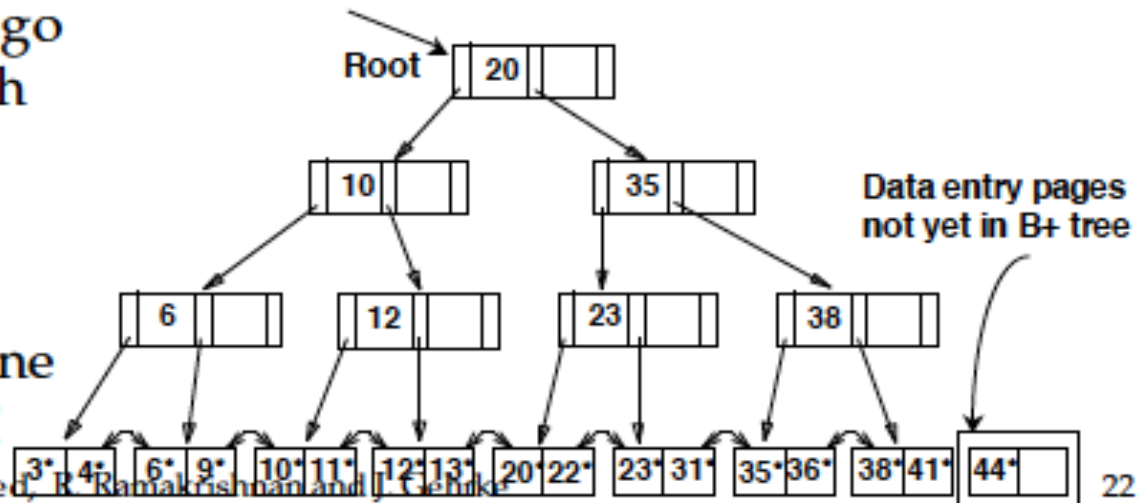
# Bulk Loading (Contd.)

❖ Index entries for leaf pages always entered into right-most index page just above leaf level. When this fills up, it splits. (Split may go up right-most path to the root.)

❖ Much faster than repeated inserts, especially when one considers locking!



Root
10 | 20

6 | | | 12 | | | 23 | 35

Data entry pages not yet in B+ tree

3* 4* | 6* 9* | 10* 11* | 12* 13* | 20* 22* | 23* 31* | 35* 36* | 38* 41* | 44*

Root
20

10 | | | 35

6 | | | 12 | | | 23 | | | 38

Data entry pages not yet in B+ tree

3* 4* | 6* 9* | 10* 11* | 12* 13* | 20* 22* | 23* 31* | 35* 36* | 38* 41* | 44*

22

# Summary

- Tree-structured indexes are ideal for range-searches, also good for equality searches.

- B+ tree is a dynamic structure.
  - Inserts/deletes leave tree height-balanced; High fanout (**F**) means depth rarely more than 3 or 4.
  - Almost always better than maintaining a sorted file.
  - Typically, 67% occupancy on average.
  - If data entries are data records, splits can change rids!

- Most widely used index in database management systems because of its versatility.  One of the most optimized components of a DBMS.

# More…

- Hash-based Indexes
  - Static Hashing
  - Extendible Hashing
  - Linear Hashing
- Grid-files
- R-Trees
- etc…
- A nice animation site for B+ trees:
  https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html