# Tree Indexes
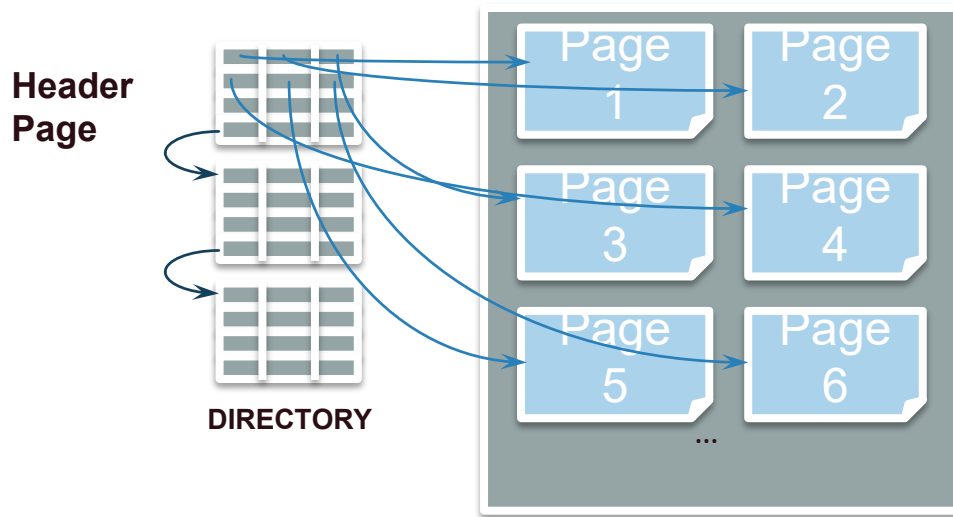
R & G - Chapter 10

# Reminder on Heap Files

- Two access APIs:
  - fetch by recordId (pageId, slotId)
  - scan (starting from some page)
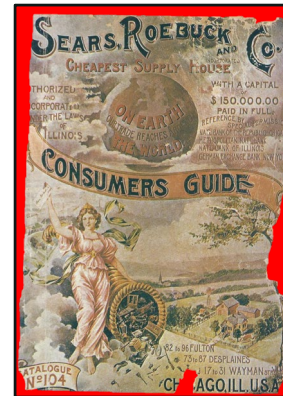
**Header Page**

**DIRECTORY**

Page 1

Page 2

Page 3

Page 4

Page 5

Page 6

...

# Wouldn't it be nice…



- …if we could look things up by value?

- Toward a Declarative access API

- But … efficiency?
  "If you don't find it in the index, look very carefully through the entire catalog. "
  —Sears, Roebuck, and Co., Consumers' Guide, 1897

# We've seen this before

- Data structures … in RAM:
  - Search trees (Binary, AVL, Red-Black, …)
  - Hash tables

- Needed: disk-based data structures
  - "paginated": made up of disk pages!

# Index

An **index** is data structure that enables fast **lookup** and **modification** of **data entries** by **search key**

- **Lookup**: may support many different operations
  - **Equality**, 1-d range, 2-d region, …

- **Search Key:** any subset of columns in the relation
  - Do not need to be unique
    - —e.g. (firstname) or (firstname, lastname)

# Index Part 2

An **index** is data structure that enables fast **lookup** and **modification** of **data entries** by **search key**

- **Data Entries:** items stored in the index
  - Assume for today: a pair (**k**, recordId) …
    - Pointers to records in Heap Files!
    - Easy to generalize later
- **Modification:** want to support fast insert and delete

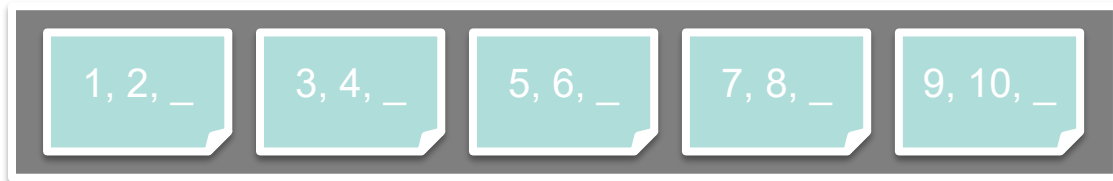Many Types of indexes exist: B+-Tree, Hash, R-Tree, GiST, ...

# Simple Idea?

Input
Heap
File

| 3, 4, 5 | 1, 2, 7 | 8, 6, 9 | 10, _, _ |

- **Step 1:** Sort heap file & leave some space

  - Pages physically stored in logical order (sequential access)
  - Do we need "next" pointers to link pages?
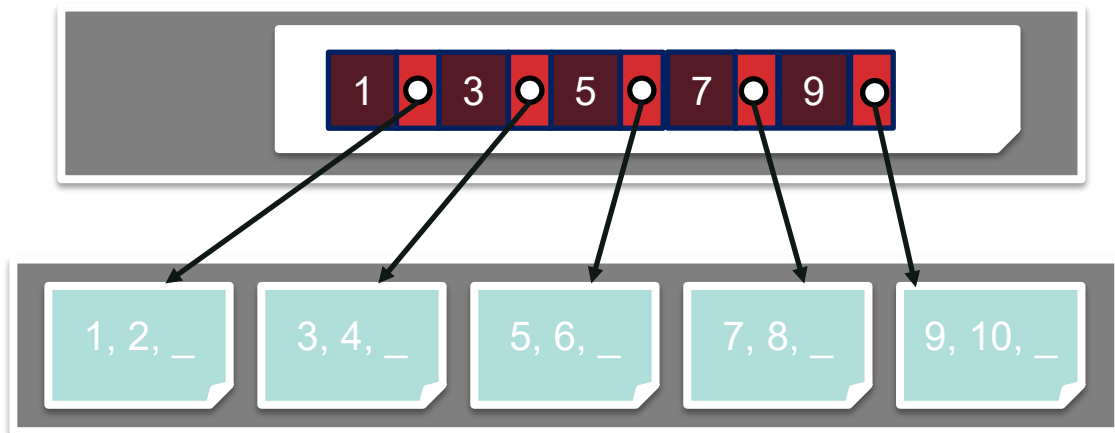    - No. Pages are physically sorted in logical order

| 1, 2, _ | 3, 4, _ | 5, 6, _ | 7, 8, _ | 9, 10, _ |

- **Step 2**: Build the index data structure over this…
  - Why not just use binary search in this heap file?
    - Fan-out of 2 → deep tree → lots of I/Os
    - Examine entire records just to read key during search
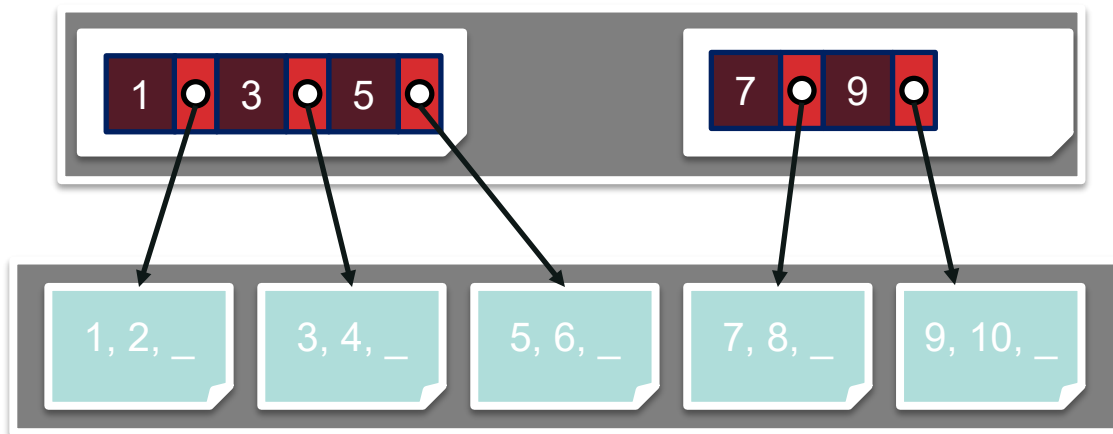
# Build a high fan-out search tree

- Start simple: *Sorted (key, page id) file*
  - No record data
  - Binary search in the key file. Better!
  - **Forgot:** Need to break across pages!

# Build a high fan-out search tree

- Start simple: *Sorted (key, page id) file*
  - No record data
  - Binary search in the key file. Better!
  - **Forgot:** Need to break across pages!
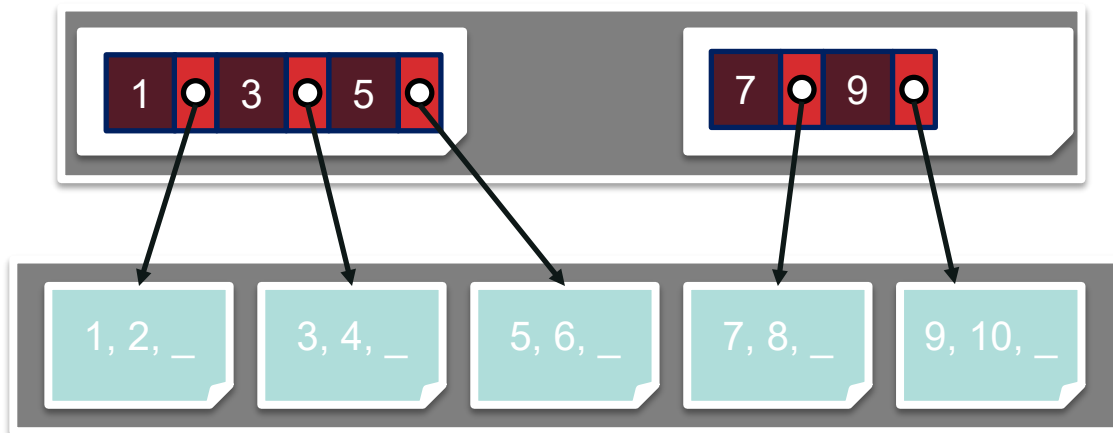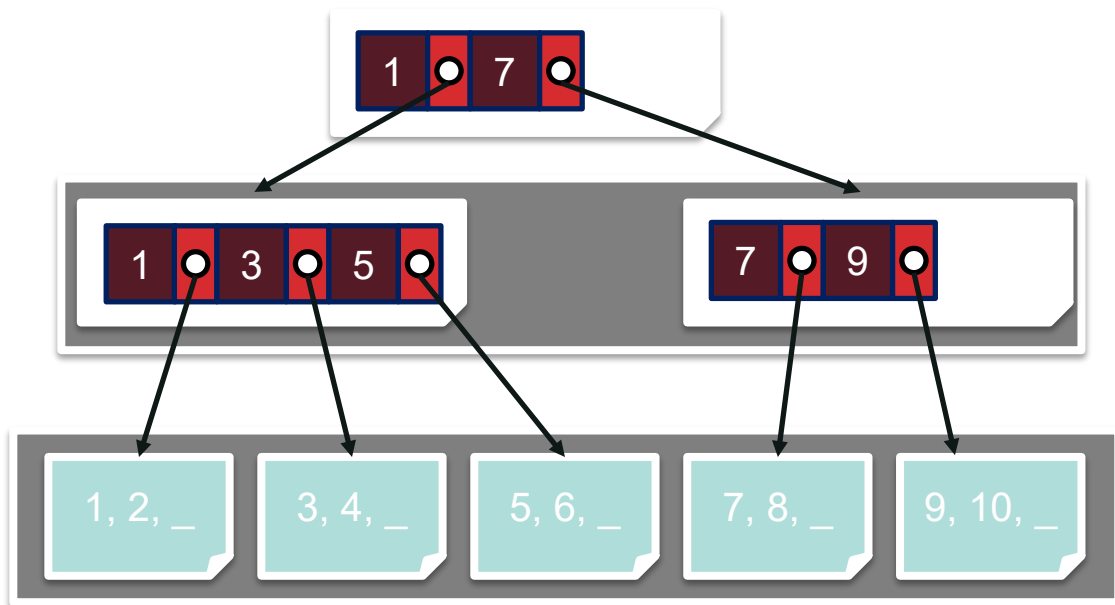
# Build a high fan-out search tree Part 2

- Start simple: *Sorted (key, page id) file*
  - No record data
  - Binary search in the key file. Better!
  - **Complexity?**

# Build a high fan-out search tree Part 3

- Start simple: *Sorted (key, page id) file*
  - No record data
  - Binary search in the key file. Better!
  - **Complexity**: Still binary search, just a constant factor smaller input
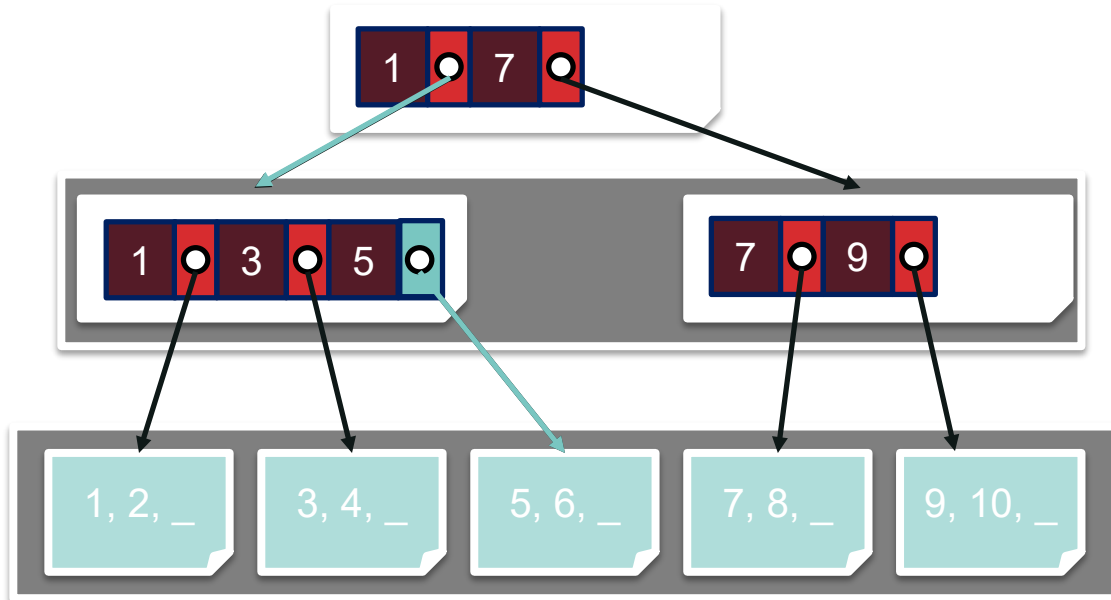
# Build a high fan-out search tree Part 4

- Recursively "index" key file
- **Key Invariant:**
  - Node […, $(K_L, P_L)$, $(K_R, P_R)$, … ] $\rightarrow$ All tuples in range $K_L <= K < K_R$ are in tree $P_L$
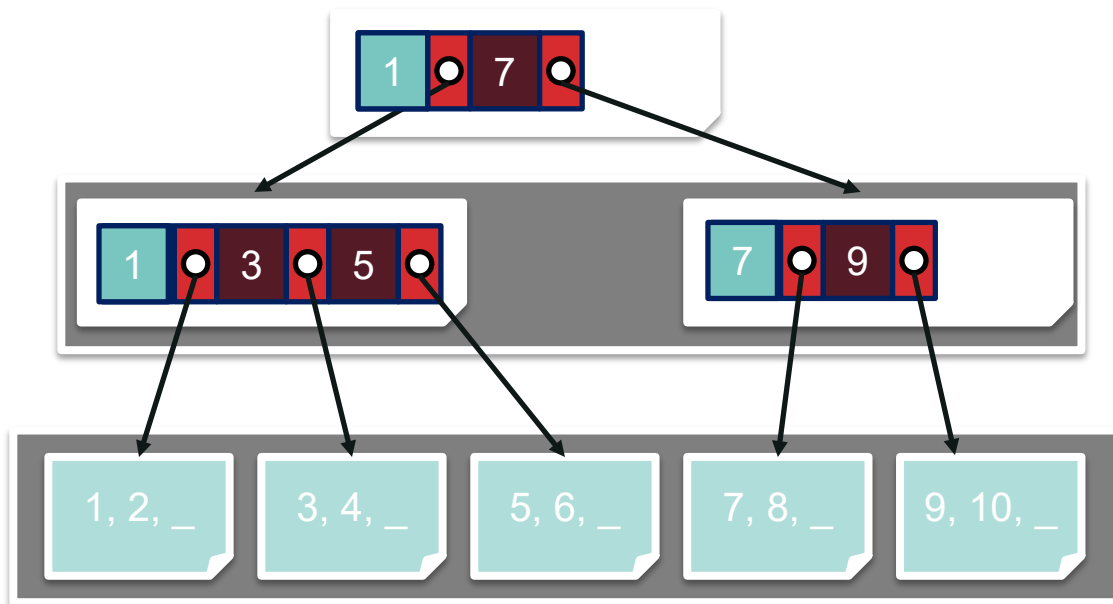
# Search a high fan-out search tree

- Searching for **5**?
  - Binary Search each node (page) starting at root
  - Follow pointers to next level of search tree
- Complexity?  O($\log_F$(#Pages))
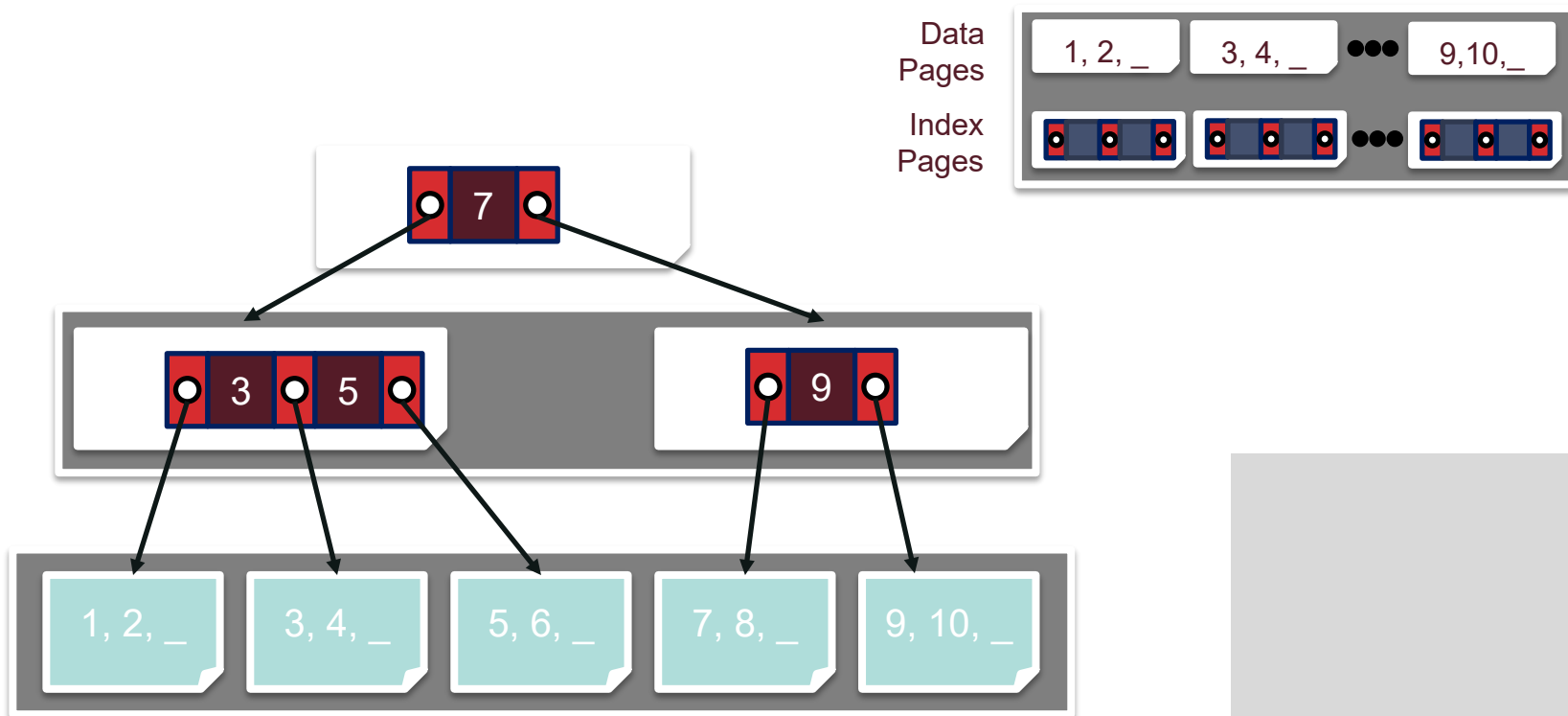
# Left Key Optimization?

- Optimization
  - Do we need the left most key?

# Build a high fan-out search tree

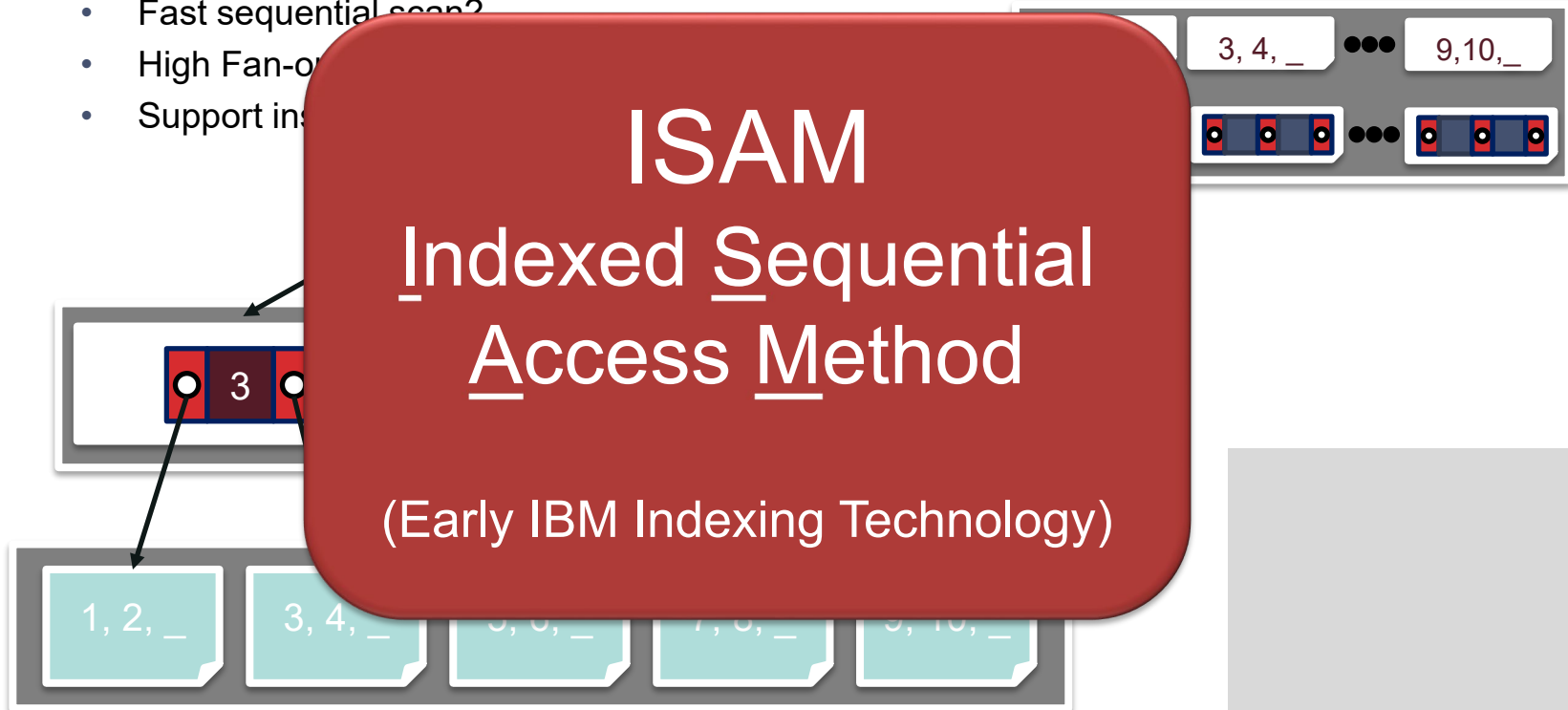- Disk Layout? All in a single file, Data Pages first.

# Status Check

- Some design goals:
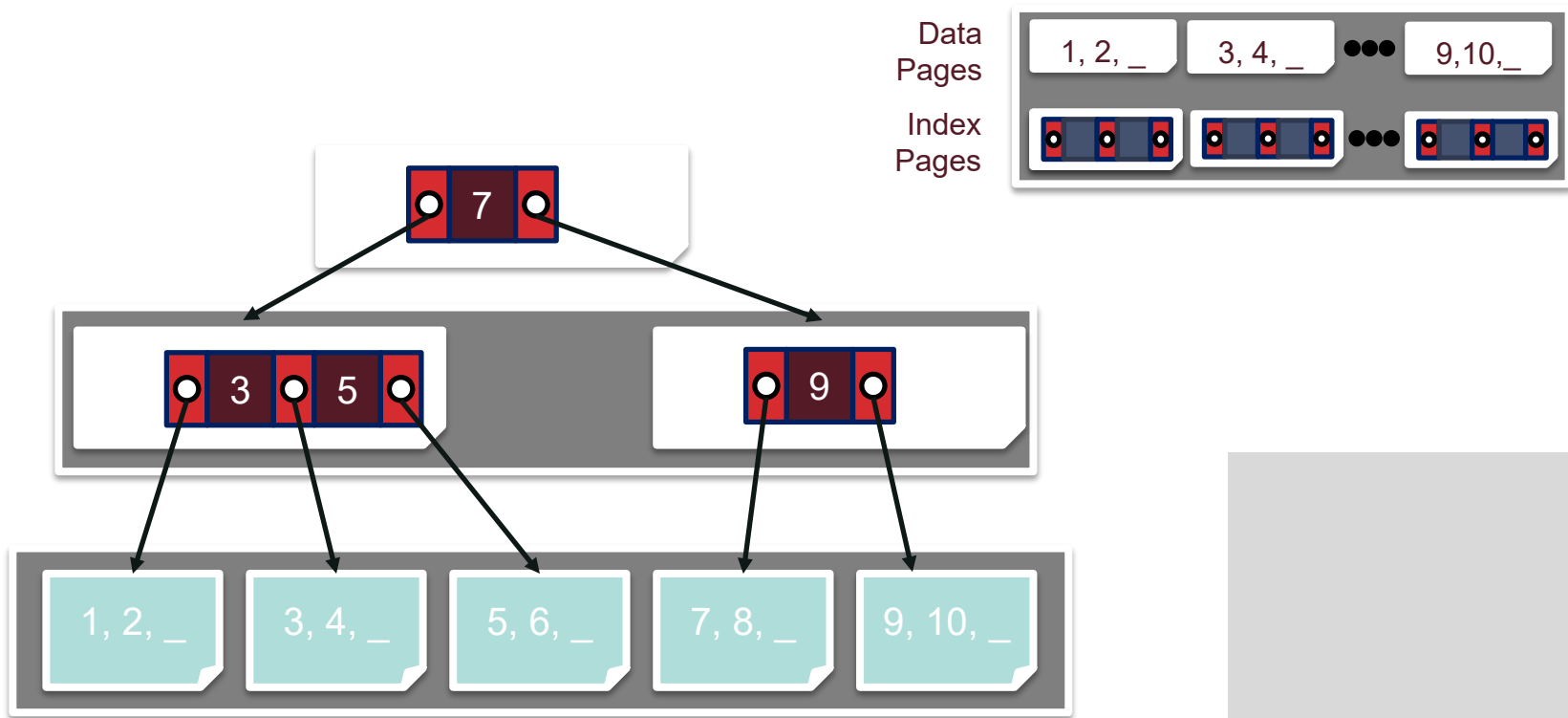  - Fast sequential scan?
  - High Fan-o~~~
  - Support ins~~~

Indexed File

3, 4, _ ●●● 9,10,_

3

1, 2, _    3, 4, _    ~~~, ~~~, _    7, 8, _    9, 10,_

**ISAM**
**Indexed Sequential**
**Access Method**

**(Early IBM Indexing Technology)**

# Insert 11, Before

Data Pages | 1, 2, _ | 3, 4, _ | ●●● | 9, 10, _

Index Pages

7

3  5

9

1, 2, _   3, 4, _   5, 6, _   7, 8, _   9, 10, _

# Insert 11, After
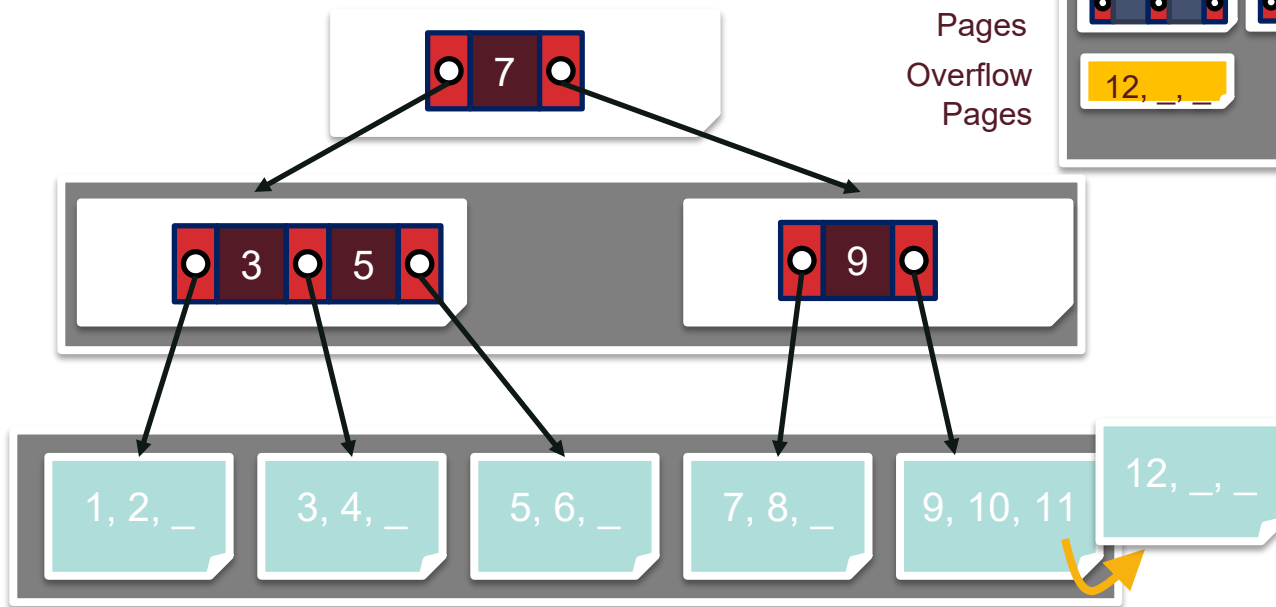
- Find location
- Place in data page
  - Re-sort page …

# Insert 12?

- Find location
- Place in data page
- Add overflow page if necessary …

# Recap: ISAM

- Data entries in sorted heap file

- High fan-out static tree index

- Fast search + good locality
  - Assuming nothing changes

- Insert into overflow pages

# A Note of Caution

- ISAM is an old-fashioned idea
  - Introduced by IBM in 1960s
  - B+ trees are usually better, as we'll see
    - Though not always (← we'll come back to this)

- But, it's a good place to start
  - Simpler than B+ tree, many of the same ideas

- Upshot
  - Don't brag about ISAM on your resume
  - Do understand ISAM, and tradeoffs with B+ trees

# B+-TREE

# Enter the B+ Tree

- Similar to ISAM
  - Same interior node structure
    - <Key, Page Ptr> pairs with same key invariant
  - Same search routine as before

- **Dynamic Tree Index**
  - Always Balanced
  - Support efficient insertion & deletion
    - Grows at root not leaves!

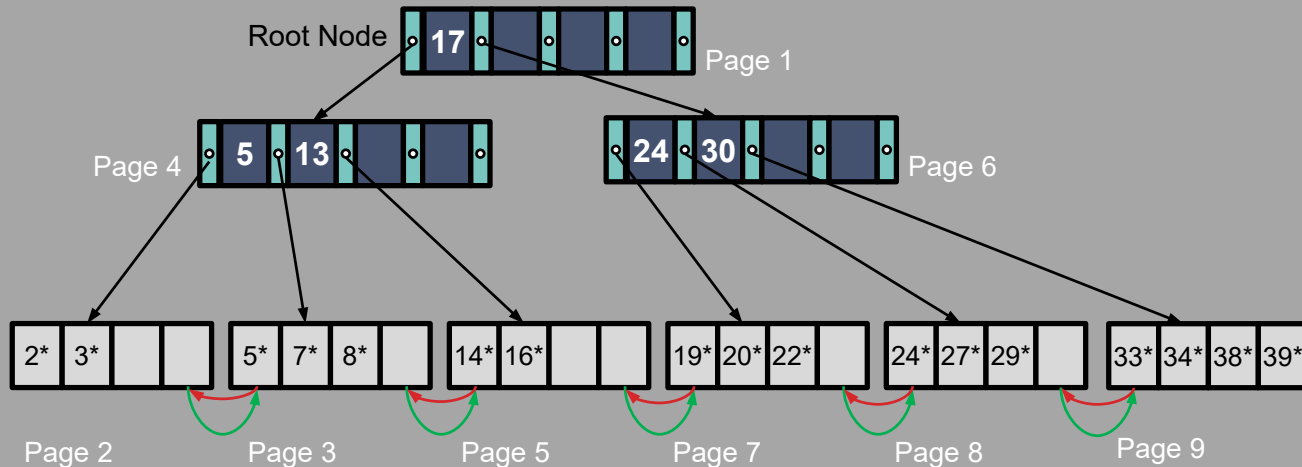- "+"? B-tree that stores data entries in leaves only

# Example of a B+ Tree



- Occupancy Invariant
  - Each interior node is at least partially full:
    - **d <= #entries <= 2d**
    - **d: order of the tree (max fan-out = 2d + 1)**
- Data pages at bottom need not be stored in logical order
  - Next and prev pointers

# Sanity Check



What is the value of d?

     2

What about the root?

     The root is special

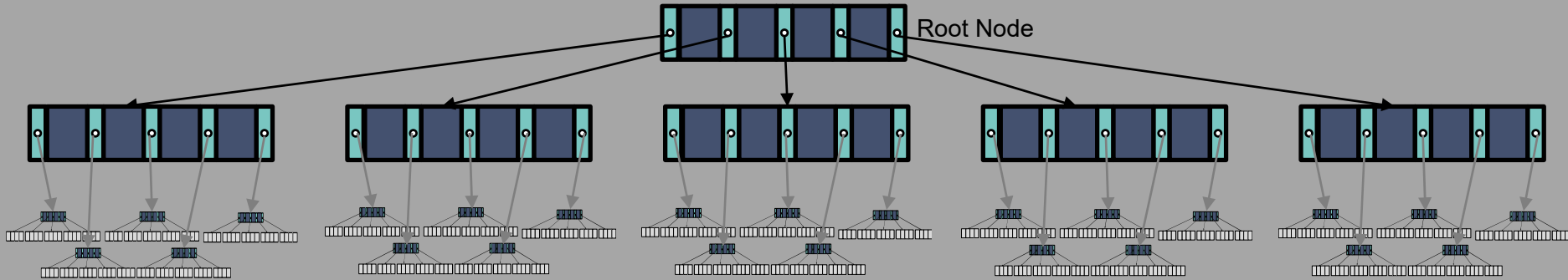Why not in sequential order?

     Data pages allocated dynamically

# B+ Trees and Scale



Root Node

- How big is a height 1 B+ tree
  - d = 2 → Fan-out?
  - Fan-out = 2d + 1 = 5
  - **Height 1:** 5 x 4 = 20 Records

# B+ Trees and Scale Part 2



- How big is a height 3 B+ tree
  - d = 2 → Fan-out?
  - Fan-out = 2d + 1 = 5
  - **Height 3:** $5^3$ x 4= 500 Records

# B+ Trees in Practice

- Typical order: 1600. Typical fill-factor: 67%.
  - average fan-out = 2144
  - (assuming 128 Kbytes pages at 40Bytes per record)

- At typical capacities
  - Height 1: $2144^2$ = **4,596,736 records**
  - Height 2: $2144^3$ = **9,855,401,984 records**

# Searching the B+ Tree



- Same as ISAM
- Find key = 27
  - Find split on each node (Binary Search)
  - Follow pointer to next node

# Searching the B+ Tree: Find 27



- Same as ISAM
- Find key = 27
  - Find split on each node (Binary Search)
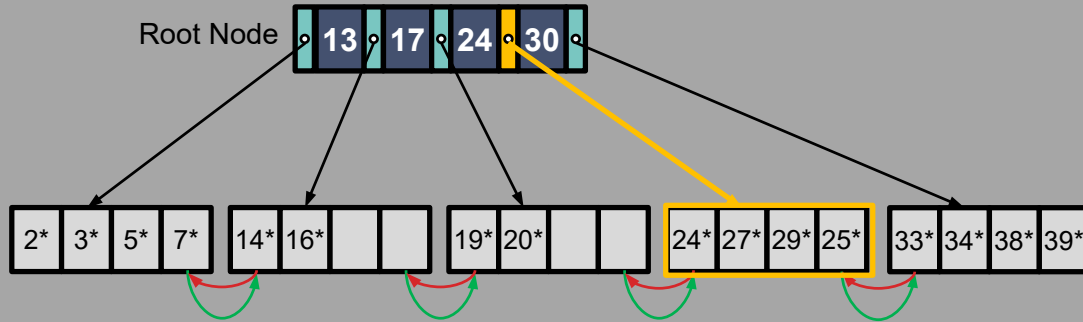  - Follow pointer to next node

# Searching the B+ Tree: Fetch Data

# Inserting 25* into a B+ Tree Part 1



- Find the correct leaf

# Inserting 25* into a B+ Tree Part 2



- Find the correct leaf
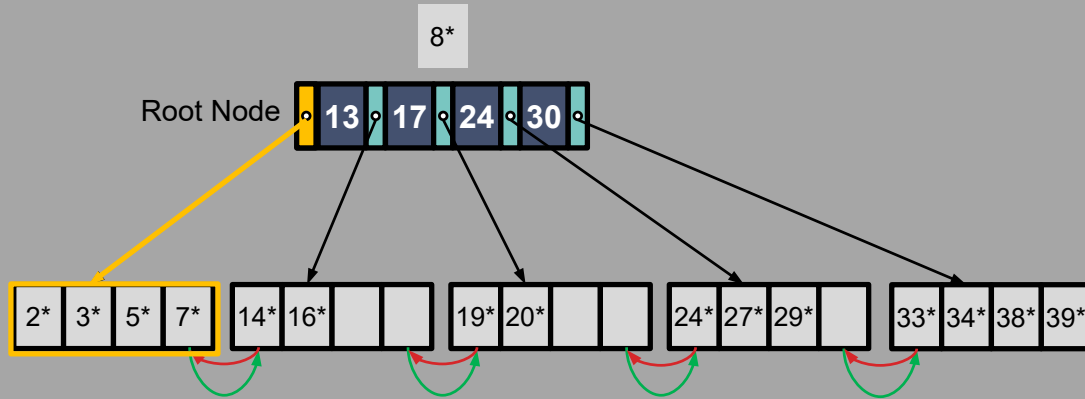- If there is room in the leaf just add the entry

# Inserting 25* into a B+ Tree Part 3



- Find the correct leaf
- If there is room in the leaf just add the entry
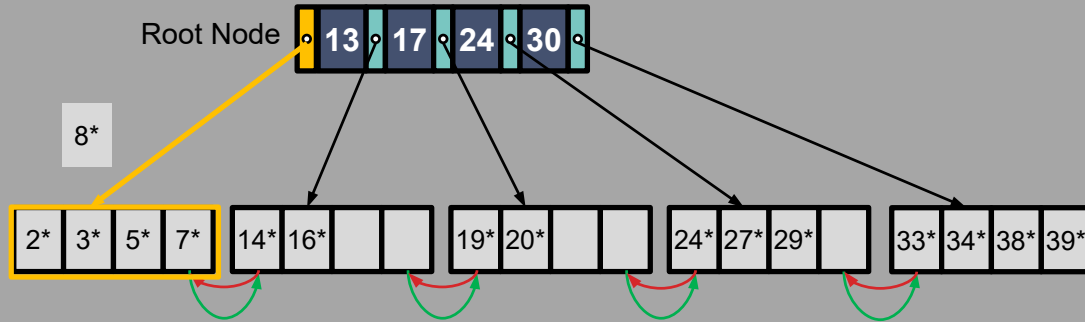  - Sort the leaf page by key

# Inserting 8* into a B+ Tree: Find Leaf
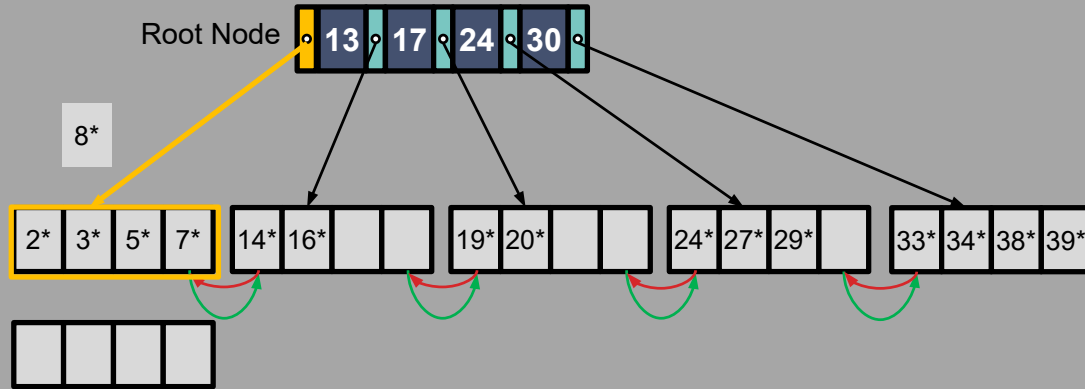


- Find the correct leaf

# Inserting 8* into a B+ Tree: Insert



- Find the correct leaf
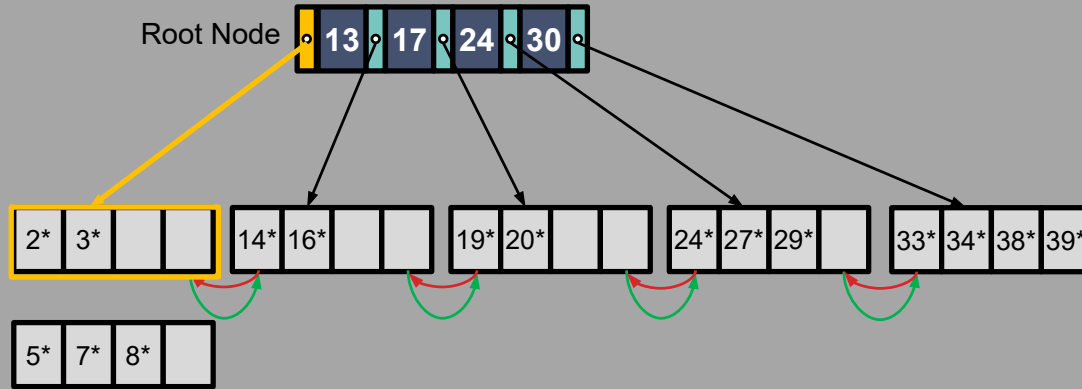  - Split leaf if there is not enough room

# Inserting 8* into a B+ Tree: Split Leaf



- Find the correct leaf
  - Split leaf if there is not enough room
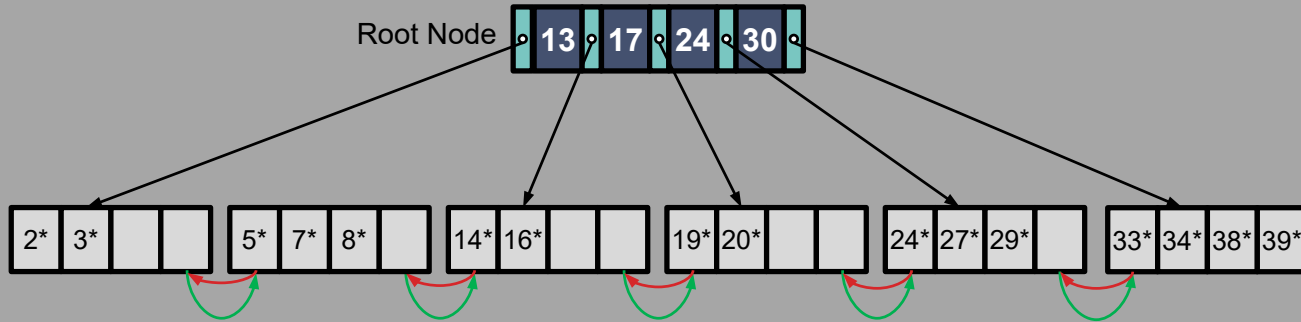  - Redistribute entries evenly

# Inserting 8* into a B+ Tree: Split Leaf, cont



- Find the correct leaf
  - Split leaf if there is not enough room
  - Redistribute entries evenly
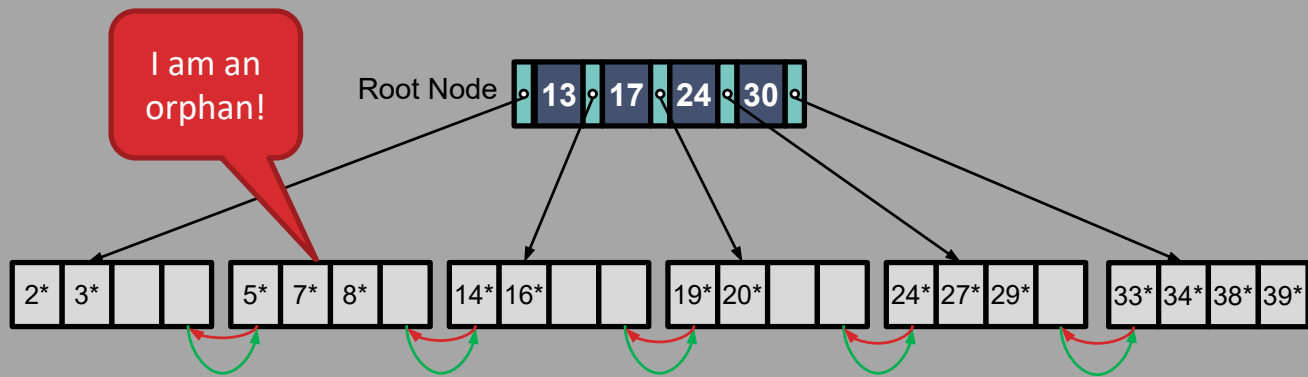  - Fix next/prev pointers

# Inserting 8* into a B+ Tree: Fix Pointers



- Find the correct leaf
  - Split leaf if there is not enough room
  - Redistribute entries evenly
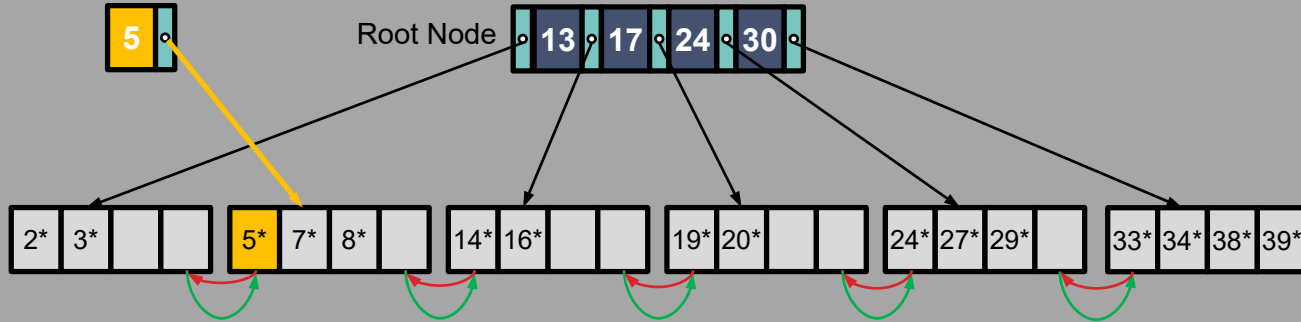  - Fix next/prev pointers

# Inserting 8* into a B+ Tree: Mid-Flight
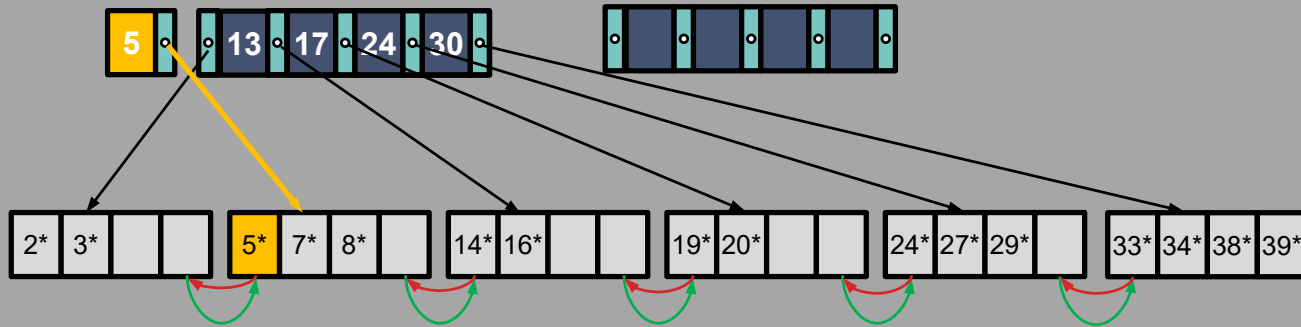


- Something is still wrong!

# Inserting 8* into a B+ Tree: Copy Middle Key



- **Copy up from leaf** the middle key
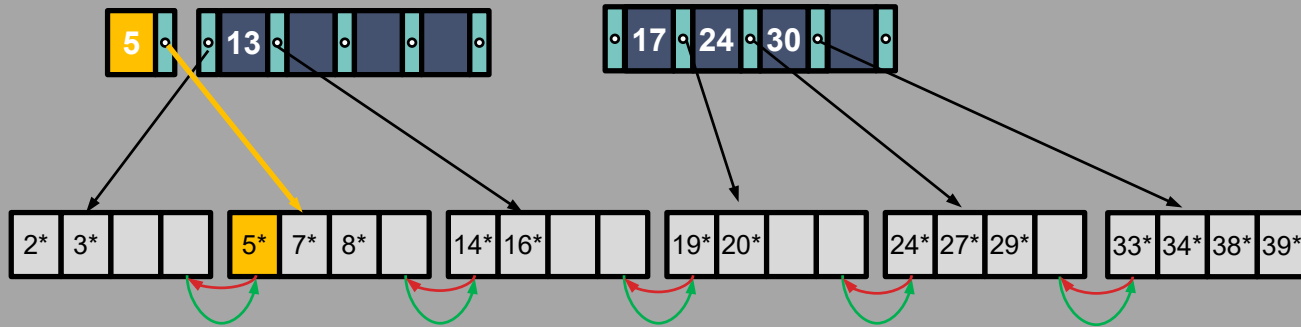- No room in parent? Recursively split index nodes

# Inserting 8* into a B+ Tree: Split Parent, Part 1



- **Copy up from leaf** the middle key
- No room in parent? Recursively split index nodes
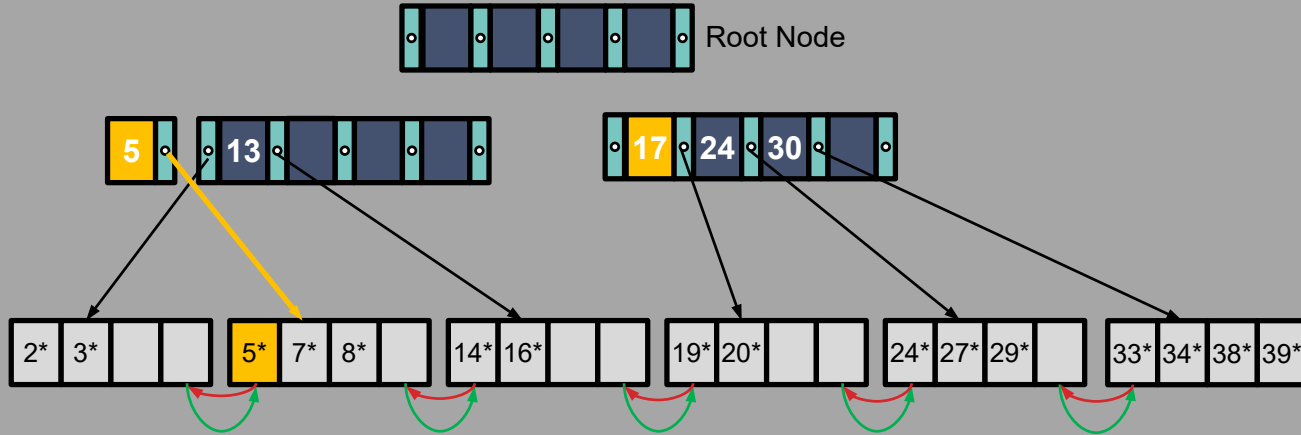  - Redistribute the rightmost d keys

# Inserting 8* into a B+ Tree: Split Parent, Part 2



- **Copy up from leaf** the middle key
- No room in parent? Recursively split index nodes
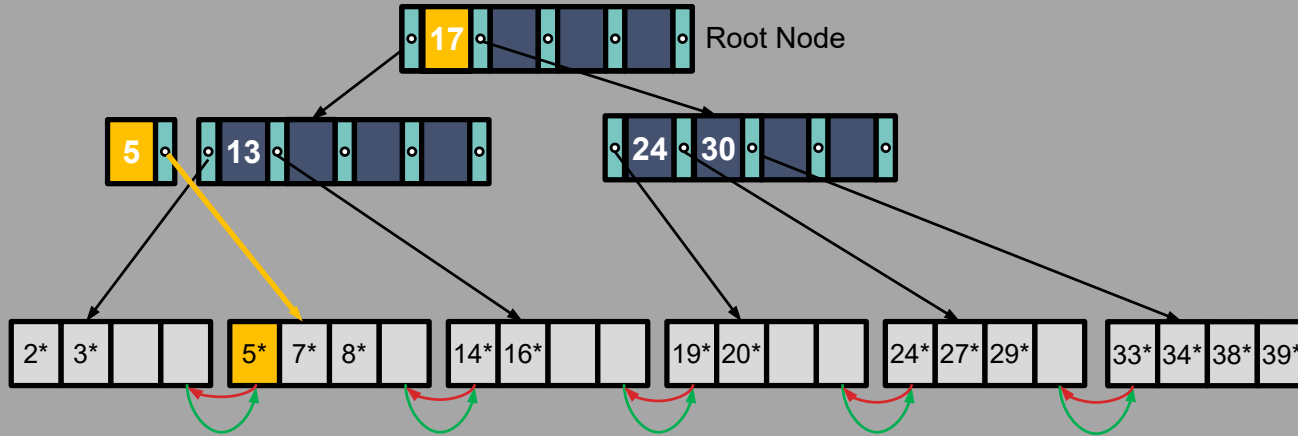  - Redistribute the rightmost d keys

# Inserting 8* into a B+ Tree: Root Grows Up



- **Push up from interior node** the middle key
  - Now the last key on left
- No room in parent? Recursively split index nodes
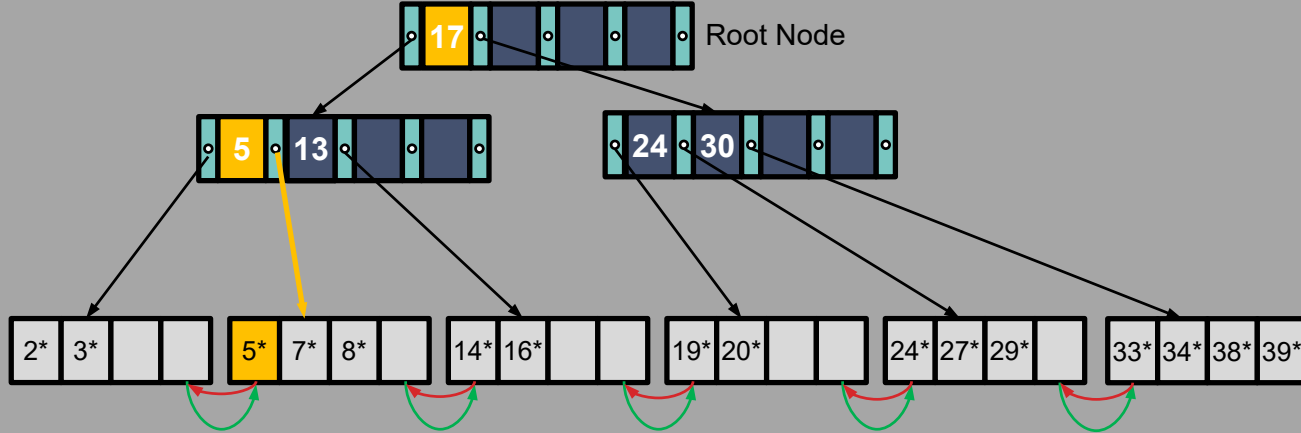  - Redistribute the rightmost d keys

# Inserting 8* into a B+ Tree: Root Grows Up, Pt 2



- Recursively split index nodes
  - Redistribute right d keys
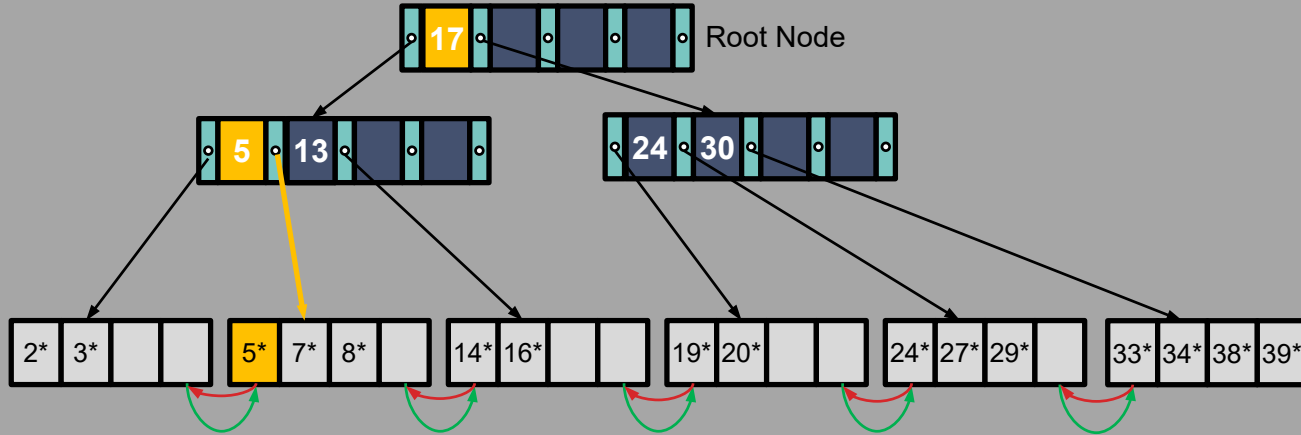  - **Push** up middle key

# Inserting 8* into a B+ Tree: Root Grows Up, Pt 3



- Recursively split index nodes
  - Redistribute right d keys
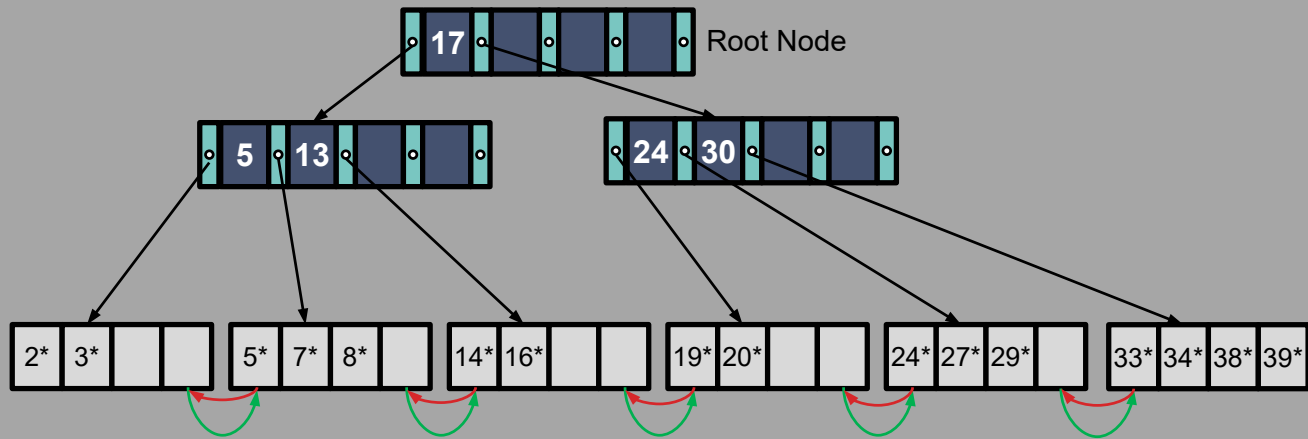  - **Push** up middle key

# Copy up vs Push up!



- Notice:
  - The **leaf** entry (5) was **copied** up
  - The **index** entry (17) was **pushed** up

# Inserting 8* into a B+ Tree: Final



- Check invariants
- **Key Invariant:**
  - Node[…, $(K_L, P_L)$, …] ➔
    $K_L <= K$ for all $K$ in $P_L$ Sub-tree
- **Occupancy Invariant:**
  - d <= # entries <= 2d