

CS 250: FOUNDATIONS OF COMPUTER SYSTEMS

[DATA STRUCTURES FOR STORAGE]

Indexes to the rescue!

Have you a surfeit of on-disk data?

Needing searches by the sweat of your brow

Maintain indexes alongside your data

Separate and apart

Update during ingestion or writes

A friend to consult when you search

Their raison d'etre

To speed up what you seek

Without many a disk seek

The search load lightened

SHRIDEEP PALLICKARA

Computer Science

Colorado State University

COMPUTER SCIENCE DEPARTMENT



COLORADO STATE UNIVERSITY

1

Frequently asked questions from the previous class survey

- Is BST used in a lot of places?
 - What can we do about unbalanced BSTs?



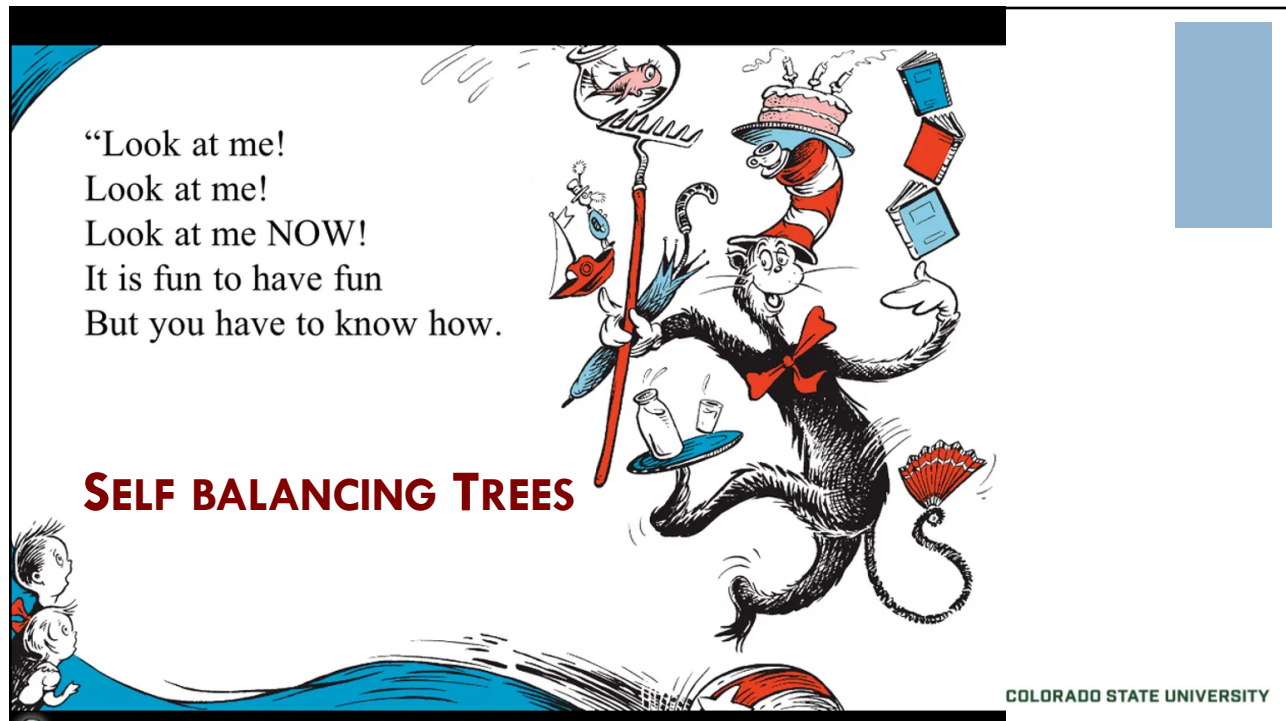
COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.2

2



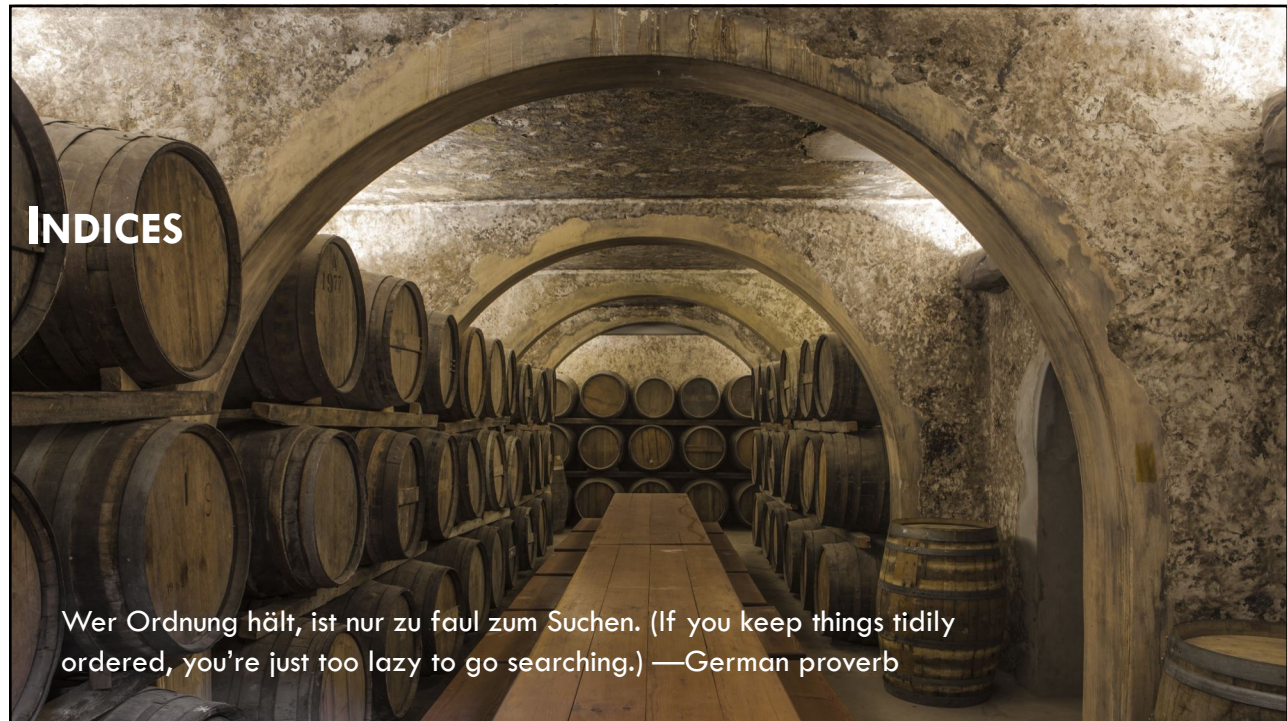
3

Topics covered in this lecture

- B-Trees

COLORADO STATE UNIVERSITY Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT DATA STRUCTURES FOR STORAGE L26.4

4



5

Let's look at the simplest SQL query out there

- We are looking for software engineers between 18-25
 - ▣ With a remuneration between 85-105K

```
SELECT *  
FROM people  
WHERE age BETWEEN 18 AND 25  
      AND salary BETWEEN 85000 AND 105000  
      AND job_title = 'Software Engineer';
```



6

```
CREATE TABLE people (  
  id INT PRIMARY KEY,  
  name VARCHAR(50),  
  age INT,  
  salary INT,  
  job_title VARCHAR(50)  
);
```

```
INSERT INTO people (id, name, age, salary, job_title) VALUES  
(1, 'Ava Chen', 19, 87000, 'Software Engineer'),  
(2, 'Liam Patel', 22, 92000, 'Software Engineer'),  
(3, 'Noah Garcia', 24, 104000, 'Software Engineer'),  
(4, 'Emma Johnson', 25, 105000, 'Software Engineer'),  
(5, 'Olivia Brown', 18, 85000, 'Software Engineer'),  
(6, 'Mason Lee', 21, 83000, 'Software Engineer'),  
(7, 'Sophia Davis', 26, 98000, 'Software Engineer'),  
(8, 'Ethan Wilson', 23, 110000, 'Software Engineer'),  
(9, 'Isabella Moore', 20, 99000, 'Data Analyst'),  
(10, 'James Taylor', 17, 90000, 'Intern'),  
(11, 'Mia Anderson', 25, 88000, 'QA Engineer'),  
(12, 'Benjamin Thomas', 24, 106000, 'Software Engineer');
```

```
SELECT *  
FROM people  
WHERE age BETWEEN 18 AND 25  
AND salary BETWEEN 85000 AND 105000  
AND job_title = 'Software Engineer';
```

```
(1, 'Ava Chen', 19, 87000, 'Software Engineer')  
(2, 'Liam Patel', 22, 92000, 'Software Engineer')  
(3, 'Noah Garcia', 24, 104000, 'Software Engineer')  
(4, 'Emma Johnson', 25, 105000, 'Software Engineer')  
(5, 'Olivia Brown', 18, 85000, 'Software Engineer')
```

7

Databases and indices

- In order to efficiently find the value for a particular key in the database, we *need* a different data structure: an **index**
- Key enabling idea
 - Keep some additional *metadata* on the side
 - **Index acts as a signpost** and helps you to locate the data you want
- If you want to search records in *several different* ways?
 - You may need several different **indexes on different parts** of the data



8

What is an index?

- An index is an additional structure that is **derived** from the primary data
- Makes searches very efficient ...



Nothing in life is free, especially fast reads

- Many databases allow you to add and remove indexes
 - ▣ This doesn't affect the contents of the database; it only affects the **performance of queries**
- Maintaining additional structures **incurs overheads during writes**
 - ▣ For writes, it's hard to beat the performance of simply appending to a file, because that's the simplest possible write operation
- Any kind of index usually slows down writes
 - ▣ Because the index also needs to be *updated every time data is written*



Let's look at this concept a little more ... with “physical” books

- The chapters and pages are the **primary data**
 - The index at the back is the additional structure *derived* from that data
- The index does not change the actual book content
- It just helps you find things faster:
 - “buffer cache” → page 117
 - “binary search” → pages 52, 89, 201
 - “B-Tree” → pages 243–247
- But now imagine the author *adds* a new chapter or *rewrites* several pages?
 - The book's content may still be fine, but the index has to be updated too
 - That is the **write overhead!**



Trade-off in storage systems

- Well-chosen indexes *speed up read queries*
 - But every index *slows down writes*
- For this reason, databases **don't usually index everything** by default
 - Require you to *choose indexes manually*
 - Using your knowledge of the application's typical query patterns
 - You can then choose the indexes that give your application the greatest benefit, without introducing more overhead than necessary



So a natural question is ...

- If indexes cost us extra writes, what kind of index is worth the trouble?
 - ▣ Let's look at the next ...



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.13

13

Boeing 747

- First flew in 1969
- Top speed is Mach 0.85
- Has carried 3.8 Billion people
- Each plane has
 - ▣ Over 140 miles of wiring
 - > 3,500 pounds
 - ▣ Has over 6 million parts



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.14

14

The most widely used indexing structure?

- The **B-tree**
- Introduced in 1970 and called “ubiquitous” less than 10 years later
 - Inventors: Rudolf Bayer & Edward M. McCreight while @ Boeing Research Labs
 - “B” was not defined: Could be for “balanced”, “broad”, “Boeing”?
- B-trees have stood the test of time very well
- They **remain the standard index implementation** in *almost all* relational databases, and many nonrelational databases use them too



Nodes in B-Trees

- Are usually *also* referred to as **pages**
- Very closely aligned with block-sizes on storage devices



B-Tree Nodes (or pages)

- The node contains several keys and references to child nodes
- Each child node is responsible for a **continuous range of keys**
 - ▣ The keys indicate where the boundaries between those ranges lie
- Most databases can fit into a B-tree that is **three or four levels deep**
 - ▣ So, you don't need to follow many references to find the page you are looking for
 - ▣ A four-level tree of 4 KB pages with a branching factor (or fanout) of 500 can store up to 250 TB



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.17

17

B-TREES



We are braver than a bee, and a... longer than a tree...
Winnie the Pooh

18

How data is read from stable storage

- HDDs and SSDs address **blocks** rather than individual bytes
- Most operating systems have a block device abstraction
- When we're reading a single word from an HDD or an SSD?
 - ▣ The whole block containing it is read



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.19

19

Primary limitation and design consideration for building efficient on-disk structures

- The **cost of disk access** itself
- The smallest unit of disk operation is a **block**
- To follow a pointer to the specific location within the block, we have to **fetch an entire block**



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.20

20

If we must always read a block?

- Why don't we **change the layout of the data structure** to take advantage of it?
- Creating long dependency chains in on-disk structures greatly increases code and structure complexity
 - Much better to keep the number of pointers and their spans to a minimum
- This is why B-Tree nodes are **page-sized**:
 - Disks fetch blocks, not individual keys, so we want each trip to the disk to bring back a lot of useful search information



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.21

21

On-disk structures optimize for target storage specifics and optimize for fewer disk accesses

- Improving locality
- Optimizing the internal representation of the structure
- Reducing the number of out-of-page pointers



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.22

22

B-Trees combine these aforementioned ideas

- Account for storage characteristics (esp. the block construct)
- Increase node fanout
- Reduce tree height
- Reduce the number of node pointers
- Reduce frequency of balancing operations



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.23

23

B-Tree analogy: How to find one thing without developing a personal relationship with every book?

- Setting: vast catalog room in the library
- You first have to pick the correct cabinet
- Then the correct shelf in that cabinet
- Then the correct drawer on the shelf, and
- Then browse through the cards in the drawer to find the one you're searching for

- Similarly, a B-Tree **builds a hierarchy** that helps to *navigate* and *locate* the searched items quickly



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.24

24

Depicting nodes in BSTs versus B-Trees [1/2]

- In most of the literature, **binary tree nodes are drawn as circles**
 - Each node is responsible for just one key and splits the range into two parts
 - This level of detail is sufficient and intuitive
- B-Tree nodes are often drawn as **rectangles**
 - Pointer blocks are also shown explicitly to highlight the relationship between child nodes and separator keys



COLORADO STATE UNIVERSITY

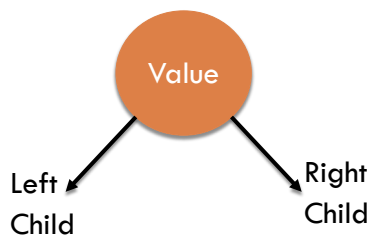
Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

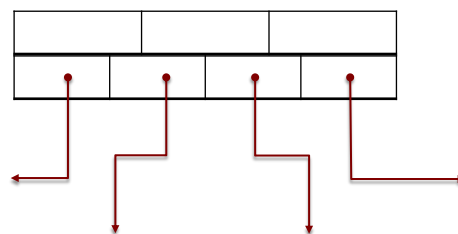
L26.25

25

Depicting nodes in BSTs versus B-Trees [2/2]



BST Node



B-Tree Node (or Page)



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.26

26

Keys inside the B-Tree nodes

- B-Trees are **sorted**
 - ▣ Keys inside the B-Tree nodes are stored in order
 - ▣ We can use an algorithm like *binary search* to locate a searched key
- This also implies that lookups in B-Trees have *logarithmic complexity*



Using B-Trees, we can efficiently execute both point and range queries

- **Point queries** locate a single item
 - ▣ Expressed by the equality (=) predicate in most query languages
- **Range queries** are used to query multiple data items in order
 - ▣ Expressed by comparison (<, >, ≤, and ≥) predicates



The B-Tree Hierarchy comprises multiple nodes

- Each node holds up to N keys
 - ▣ And N + 1 pointers to the child nodes
- Nodes are logically grouped into three groups:
 - ▣ **Root** node, which is the top of the tree
 - ▣ **Leaf** nodes: Bottommost layer nodes that have no child nodes
 - ▣ **Internal** nodes These are all other nodes with leaves
 - ▣ There is usually more than one level of internal nodes



COLORADO STATE UNIVERSITY

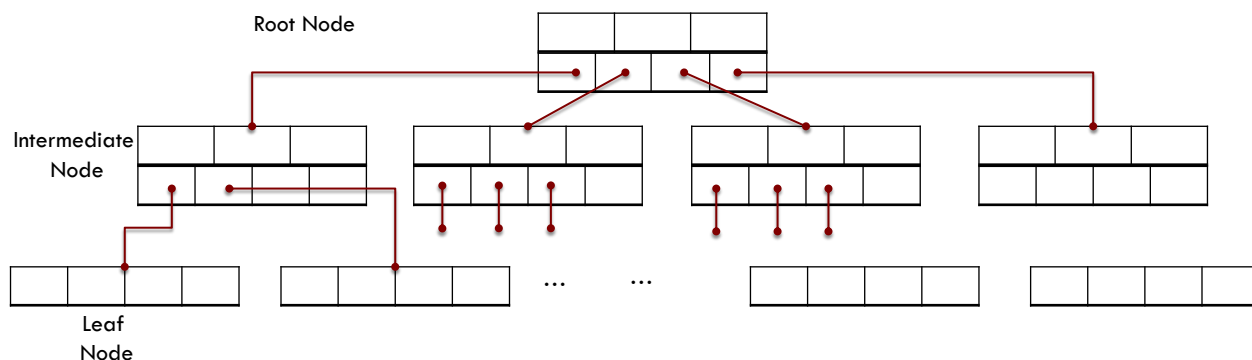
Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.29

29

B-Tree node hierarchy



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.30

30

Nomenclature

- Since B-Trees are a page organization technique
 - i.e., they are used to organize and navigate fixed-size pages
 - We often use terms node and page interchangeably
- The relation between the **node capacity** and the number of keys it actually holds is called **occupancy**



B-Trees are characterized by their fanout

- Fanout refers to the number of keys stored in each node
- **Higher fanout** helps to:
 - **Amortize** the cost of **structural changes** required to keep tree balanced
 - **Reduces the number of seeks** by storing keys and pointers to child nodes in a single block or multiple consecutive blocks
- Balancing operations (namely, splits and merges) are triggered when the nodes are full or nearly empty



Another important point ...

- High fanout helps only if each node clearly **carves** the key space into ranges
 - **Separator keys** are what perform that carving!



Separator Keys

- Keys stored in B-Tree nodes are called index entries, separator keys, or divider cells
- **Split the tree into subtrees** (also called branches or subranges), holding corresponding key ranges
 - Keys are stored in sorted order to allow binary search
- A subtree is found by locating a key and following a corresponding pointer from the higher-level to the lower-level

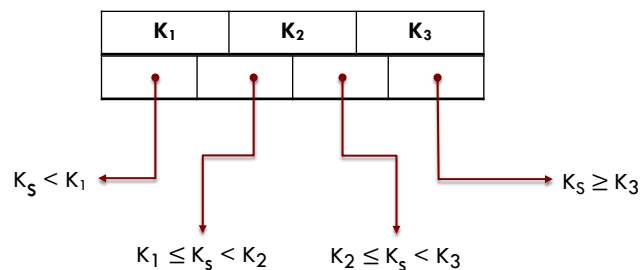


About pointers in a node

- The first pointer in the node?
 - ▣ Points to the subtree holding items less than the first key
- The last pointer in the node?
 - ▣ Points to the subtree holding items greater than or equal to the last key
- Other pointers?
 - ▣ Reference subtrees **between the two keys**: $K_{i-1} \leq K_s < K_i$, where K is a set of keys, and K_s is a key that belongs to the subtree.



Separator keys splitting a tree into subtrees



B-tree construction

- Rather than being built from top-to-bottom (as in BSTs), B-Trees are from **bottom-to-top**
- The number of leaf nodes grows, which increases the number of internal nodes and tree height
- B-Trees **reserve extra space** inside nodes for future insertions and updates
 - ▣ Tree storage utilization can get as low as 50%, but is usually much higher
- Higher occupancy does not influence B-Tree performance negatively



B-Tree lookup complexity can be viewed from two standpoints

- The number of block transfers
- The number of comparisons done during the lookup



B-Tree lookup complexity: Number of block transfers

- In terms of number of transfers, the logarithm base is N (number of keys per node)
 - There are N times more nodes on each new level
 - Following a child pointer reduces the search space by the factor of N
 - During lookup, at most $\log_N M$ pages are addressed to find a target key
 - M is the total number of items in the B-Tree
 - The number of child pointers that have to be followed on the root-to-leaf pass is also equal to the number of levels
 - In other words, the height h of the tree



B-Tree lookup complexity: Number of block transfers

- B-Tree lookup complexity is generally referenced as $\log M$
- Logarithm base is generally not used in complexity analysis
 - Changing the base simply **adds a constant factor**
 - Multiplication by a constant factor **does not change complexity**
 - For example, given the nonzero constant factor c , $O(|c| \times n) == O(n)$



B-Tree lookup complexity: Number of comparisons

- From the perspective of number of comparisons within a node
 - ▣ The logarithm base is 2
 - ▣ Since searching a key inside each node is done using binary search
 - ▣ Every comparison halves the search space
- Complexity is $\log M$



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.41

41

Different ways to describe key and child offset counts

[1/2]

- The original paper refers to device-dependent natural number k
 - ▣ Nodes, in this case, can hold between k and $2k$ keys, but can be partially filled
 - ▣ Hold at least $k + 1$ and at most $2k + 1$ pointers to child nodes
- The root page can hold between 1 and $2k$ keys
 - ▣ Later, a number l is introduced, and it is said that any nonleaf page can have $l + 1$ keys



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.42

42

Different ways to describe key and child offset counts

[2/2]

- Other sources, describe nodes that can hold up to N separator keys and $N + 1$ pointers, with otherwise similar semantics and invariants
- Both approaches bring us to the same result
 - Differences are only used to emphasize the contents of each source
 - We stick to N for clarity



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.43

43

We now know what a node looks like ...

- Next, we walk the tree the way the database would



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.44

44



B-TREE LOOKUPS

Sleep with one eye open
Gripping your pillow tight

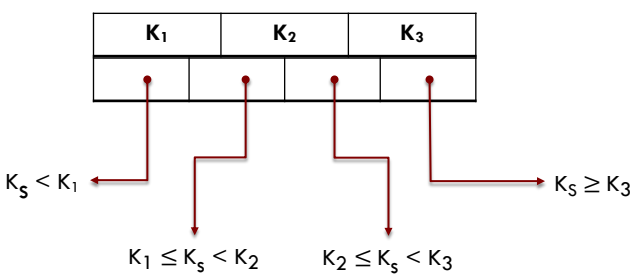
Exit light
Enter night
Take my hand
We're off to never-never land

Enter Sandman, Metallica

COMPUTER SCIENCE DEPARTMENT  COLORADO STATE UNIVERSITY


45

Separator keys splitting a tree into subtrees



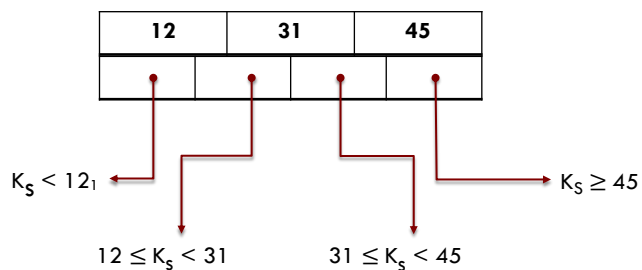
K_1	K_2	K_3
-------	-------	-------

$K_s < K_1$ $K_1 \leq K_s < K_2$ $K_2 \leq K_s < K_3$ $K_s \geq K_3$

 COLORADO STATE UNIVERSITY Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT DATA STRUCTURES FOR STORAGE L26.46

46

Separator keys splitting a tree into subtrees



COLORADO STATE UNIVERSITY

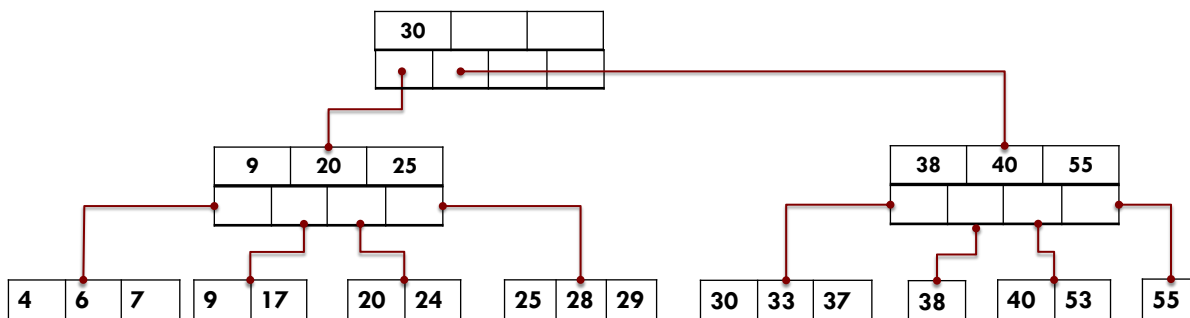
Professor: SHRIDEEP PALLICKARA
 COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.47

47

B-Tree: Example



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
 COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.48

48

B-tree Lookup Algorithm:

[1/4]

- To find an item in a B-Tree, we perform a single **traversal from root to leaf**
- The objective of this search is to find the key or its predecessor
 - ▣ Finding an exact match is used for point queries, updates, and deletions
 - ▣ Finding its predecessor is useful for range scans and inserts



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.49

49

B-tree Lookup Algorithm:

[2/4]

- Index keys split the tree into **subtrees**
 - ▣ With boundaries between two neighboring keys
- The algorithm starts from the root and performs a binary search
 - ▣ This locates a subtree
- As soon as we find the subtree?
 - ▣ Follow the pointer that corresponds to it
 - ▣ Repeat search



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.50

50

B-tree Lookup Algorithm:

[3/4]

- On each level, we get a more detailed view of the tree:
 - ▣ We start on the most coarse-grained level (the root of the tree)
 - ▣ Descend to the next level where keys represent more precise, detailed ranges
 - ▣ Until we finally reach **leaves**, where the data records are located



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.51

51

B-tree Lookup Algorithm:

[4/4]

- During the point query
 - ▣ Search completes after finding (or failing to find) the target key
- During the range scan
 - ▣ Sibling pointers are followed until the end of the range is reached or the range predicate is exhausted



COLORADO STATE UNIVERSITY

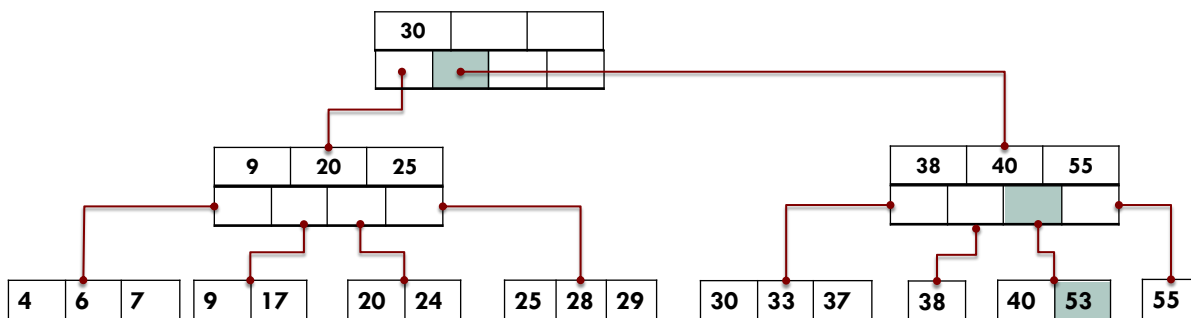
Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.52

52

B-Tree: Example: Looking for 53



COLORADO STATE UNIVERSITY

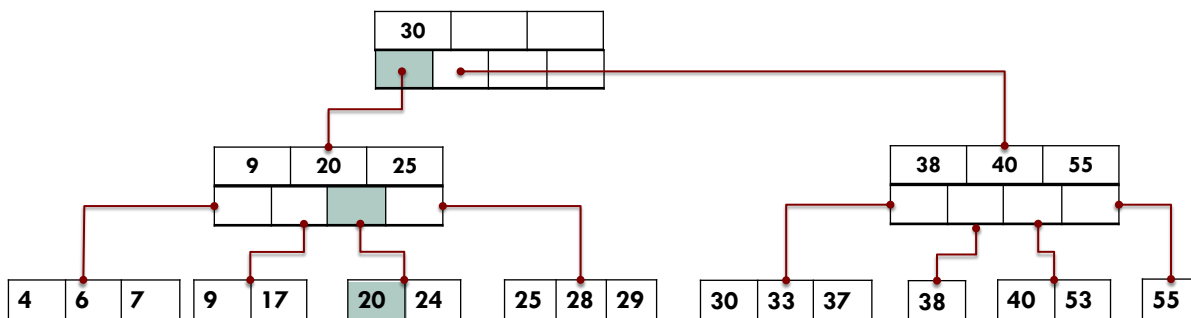
Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.53

53

B-Tree: Example: Looking for 20



COLORADO STATE UNIVERSITY

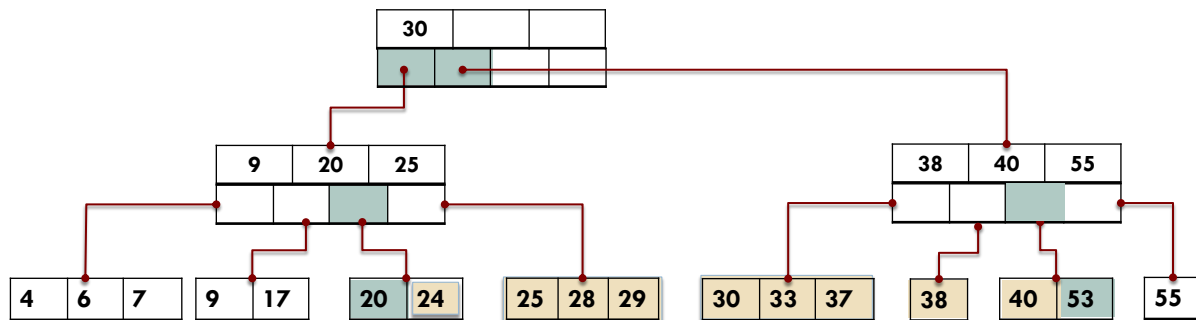
Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.54

54

B-Tree: Example: Looking for $20 < x < 53$



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.55

55

The contents of this slide-set are based on the following references

- Alex Petrov. *Database Internals*. ISBN-10/13: 1492040347/978-1492040347 O'Reilly Media. [Chapters 2,4]
- Martin Kleppmann. *Designing Data-Intensive Applications*. ISBN-10/13: 1449373321/ 978-1449373320. O'Reilly Media. [Chapter 3]



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

DATA STRUCTURES FOR STORAGE

L26.56

56