

## B+ Tree's: Primary method for managing data on disks (in RDBMS)

Objective: Understand the implications disk access on the design of data structures.

- fat fanout
- split/promote

Reading:

- Text Ch. 14
- For very helpful interactive demo, see

<https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html>

Slide thanks: many slides, Garcia-Molina, (but you may want to see his original set), <http://infolab.stanford.edu/~hector/cs245/notes.htm>

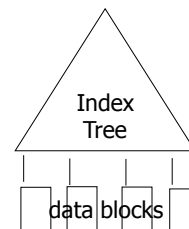
6: B+ trees

Data Engineering

1

## Data records (rows) stored in one file

Index records stored separately.



6: B+ trees

Data Engineering

2

### Sequential File

10	
20	
30	
40	
50	
60	
70	
80	
90	
100	

6: B+ trees

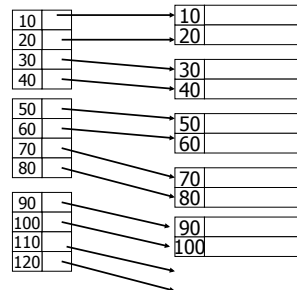
Data Engineering

3

### Dense Index

### Sequential File

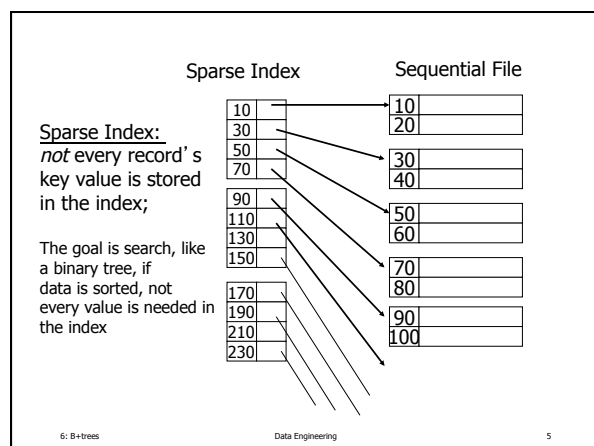
Dense Index:  
every record's key value is stored in the index



6: B+ trees

Data Engineering

4



### Sparse vs. Dense Tradeoff

- Sparse: Less index space per record can keep more of index in memory
- Dense: Can tell if any record exists without accessing file

(Later:

- sparse better for insertions
- dense needed for secondary indexes)

6: B+trees

Data Engineering

6

### Terms

- Index sequential file
- Search key ( ≠ primary key)
- Primary index (on Sequencing field)
- Secondary index
- Dense index (all Search Key values in)
- Sparse index
- Multi-level index

6: B+trees

Data Engineering

7

### The B+ Tree

- Main memory binary trees evolved to disk.
- First there were B Trees
- Then B+ Trees.
  - So superior, B Trees are gone,
    - So gone, we say B Trees and mean B+ Tree

6: B+trees

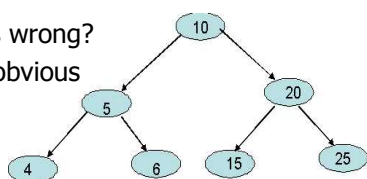
Data Engineering

8

## Naïve Use of Binary Tree on Disk

One node per disk block

- What goes wrong?
- What are obvious fixes?



A simple binary tree

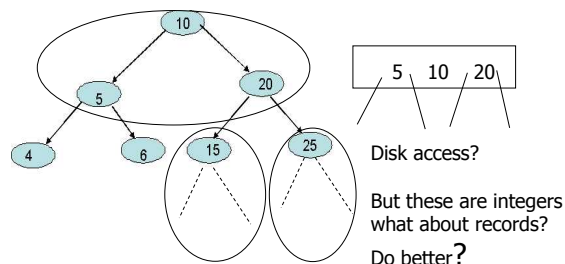
6: B+trees

Data Engineering

10

## Fatten the Nodes

(e.g. not a B+ node yet)



6: B+trees

Data Engineering

10

## In databases

- We like algorithmic complexity (why?)
- But we don't really get to throw away the constants, etc.

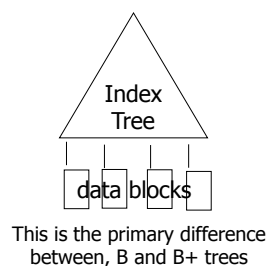
6: B+trees

Data Engineering

11

## In the index, store only search keys, not the records

- Fatter fanout - more space for *split* keys
- All the data is in the data file.
- (if dense index, all keys are in the index)

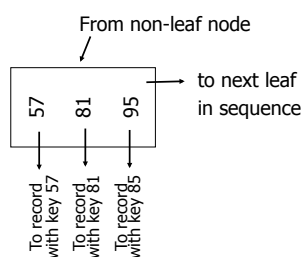


6: B+trees

Data Engineering

12

### Sample (dense) leaf node:

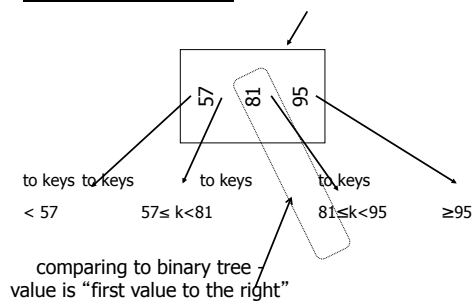


6: B+trees

Data Engineering

13

### Sample non-leaf



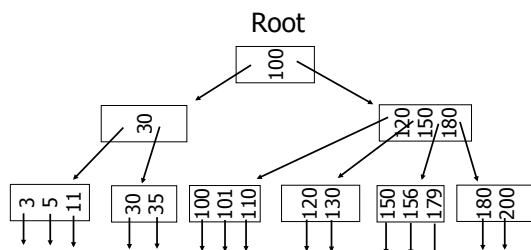
6: B+trees

Data Engineering

14

### B+ Tree Example

n=3



6: B+trees

Data Engineering

15

Size of nodes:  $\left\{ \begin{array}{l} n+1 \text{ pointers} \\ n \text{ keys} \end{array} \right.$  (fixed)

6: B+trees

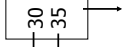
Data Engineering

16

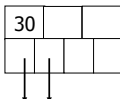
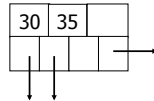
In textbook's notation

n=3

Leaf:



Non-leaf:



6: B+trees

Data Engineering

17

Don't want nodes to be too empty

- Use at least

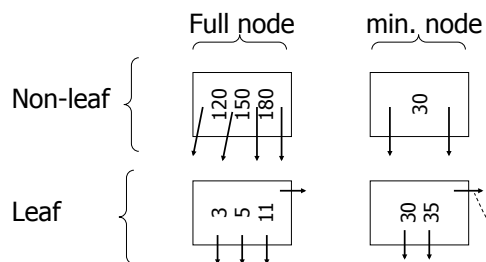
Non-leaf:  $\lceil (n+1)/2 \rceil$  pointersLeaf:  $\lfloor (n+1)/2 \rfloor$  pointers to data

6: B+trees

Data Engineering

18

n=3



6: B+trees

Data Engineering

19

**B+tree rules**

- (1) All leaves at same lowest level  
(balanced tree)
- (2) Pointers in leaves point to records  
except for "sequence pointer"

6: B+trees

Data Engineering

20

## (3) Number of pointers/keys for B+tree

	Max ptrs	Max keys	Min ptrs-data	Min keys
Non-leaf (non-root)	$n+1$	$n$	$\lceil (n+1)/2 \rceil$	$\lceil (n+1)/2 \rceil - 1$
Leaf (non-root)	$n+1$	$n$	$\lfloor (n+1)/2 \rfloor$	$\lfloor (n+1)/2 \rfloor$
Root	$n+1$	$n$	1	1

6: B+trees

Data Engineering

21

Insert into B+tree

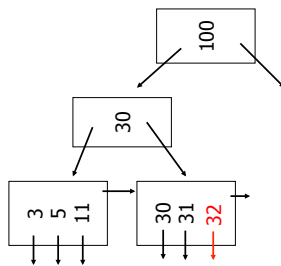
- (a) simple case
  - space available in leaf
- (b) leaf overflow
- (c) non-leaf overflow
- (d) new root

6: B+trees

Data Engineering

22

## (a) Insert key = 32

 $n=3$ 

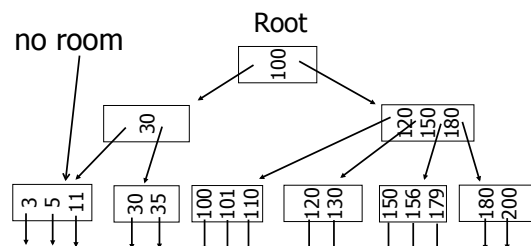
6: B+trees

Data Engineering

23

## (b) Insert key = 7

no room



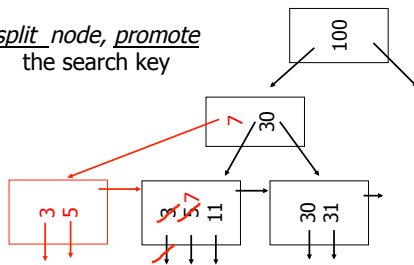
6: B+trees

Data Engineering

24

(b) Insert key = 7

n=3

split node, promote  
the search key

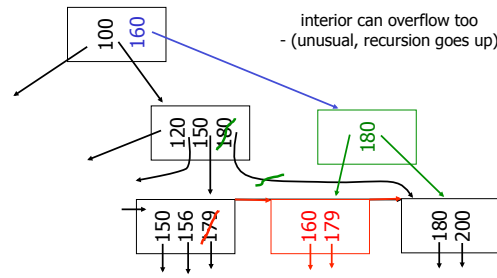
6: B+trees

Data Engineering

25

(c) Insert key = 160

n=3

interior can overflow too  
- (unusual, recursion goes up)

6: B+trees

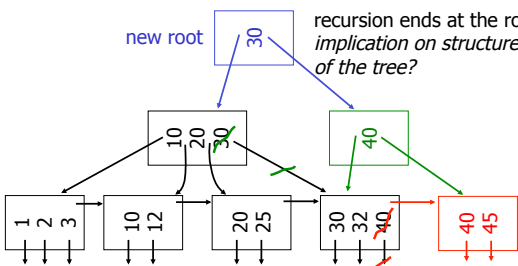
Data Engineering

26

(d) New root, insert 45

n=3

new root

recursion ends at the root  
implication on structure  
of the tree?

6: B+trees

Data Engineering

27

## Deletion from B+tree

- (a) Simple case - no example
- (b) Coalesce with neighbor (sibling)
- (c) Re-distribute keys
- (d) Cases (b) or (c) at non-leaf

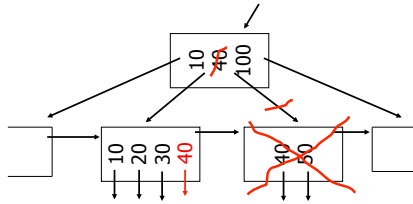
6: B+trees

Data Engineering

28

(b) Coalesce with sibling  
– Delete 50

n=4



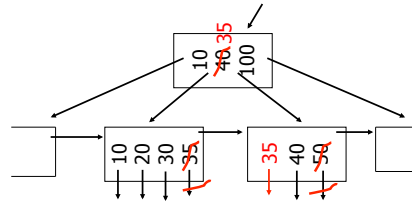
6: B+trees

Data Engineering

29

(c) Redistribute keys  
– Delete 50

n=4



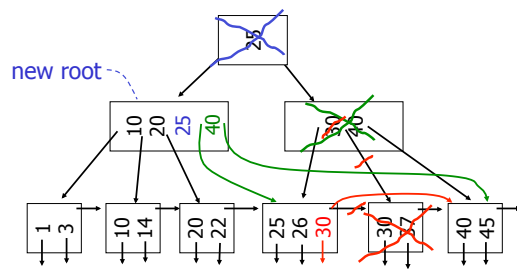
6: B+trees

Data Engineering

30

(d) Non-leaf coalesce  
– Delete 37

n=4



6: B+trees

Data Engineering

31

## B+tree deletions in practice

- Often, coalescing is not implemented
  - Overhead not worth it! (why?)

6: B+trees

Data Engineering

32



Next:

- Duplicate keys
- Deletion/Insertion
- Secondary indexes

6: B+trees

Data Engineering

33

Duplicate keys


10	
10	
10	
20	
20	
30	
30	
40	
45	

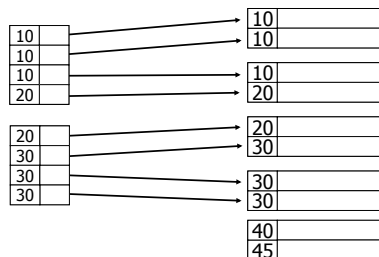
6: B+trees

Data Engineering

34

Duplicate keys

Dense index, one way to implement?



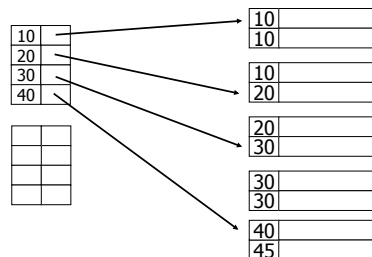
6: B+trees

Data Engineering

35

Duplicate keys

Dense index, better way?



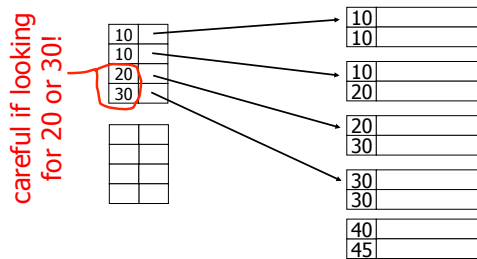
6: B+trees

Data Engineering

36

## Duplicate keys

## Sparse index, one way?



6: B+trees

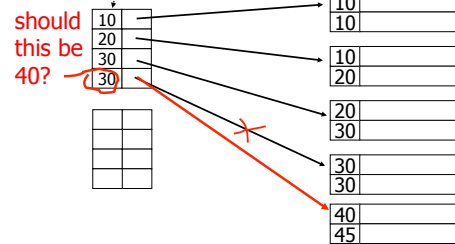
Data Engineering

37

## Duplicate keys

## Sparse index, another way?

– place first new key from block



6: B+trees

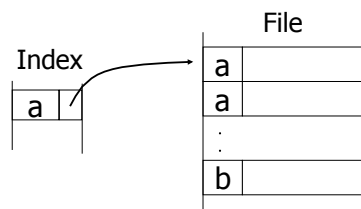
Data Engineering

38

## Summary

Duplicate values,  
primary index

- Index may point to first instance of each value only



6: B+trees

Data Engineering

39

## Primary vs. Secondary Index

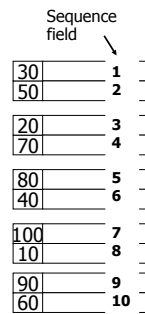
- Primary is sorted on the index key  
*and placed on disk that way*
- Secondary index,
  - does not “cluster”
  - consider (x, y)

6: B+trees

Data Engineering

40

## Secondary indexes



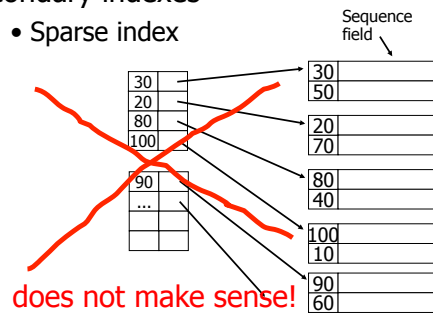
6: B+trees

Data Engineering

41

## Secondary indexes

### • Sparse index



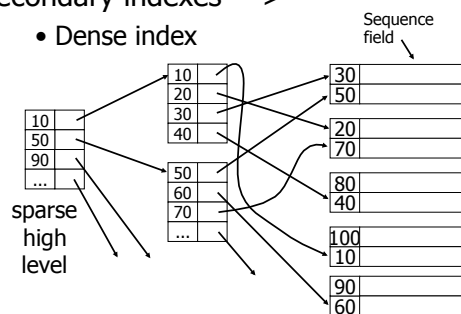
6: B+trees

Data Engineering

42

## Secondary indexes -->

### • Dense index



6: B+trees

Data Engineering

43

## With secondary indexes:

- Lowest level is dense
- Other levels are sparse

Also: Pointers are record pointers  
(not block pointers; not computed)

6: B+trees

Data Engineering

44

## Duplicate values &amp; secondary indexes

20	
10	
20	
40	
10	
40	
10	
40	
30	
40	

6: B+ trees

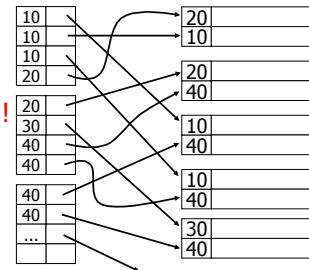
Data Engineering

45

Duplicate values & secondary indexes  
one option...

**Problem:**  
excess overhead!

- disk space
- search time



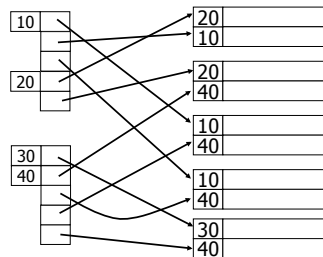
6: B+ trees

Data Engineering

46

Duplicate values & secondary indexes  
another option...

**Problem:**  
variable size  
records in  
index!

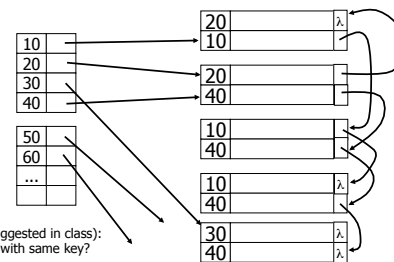


6: B+ trees

Data Engineering

47

## Duplicate values &amp; secondary indexes



Another idea (suggested in class):  
Chain records with same key?

**Problems:**

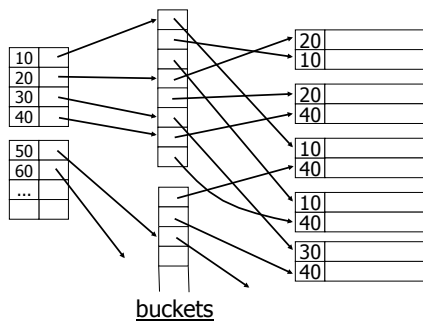
- Need to add fields to records
- Need to follow chain to know records

6: B+ trees

Data Engineering

48

### Duplicate values & secondary indexes



6: B+trees

Data Engineering

49

### Why “bucket” idea is useful

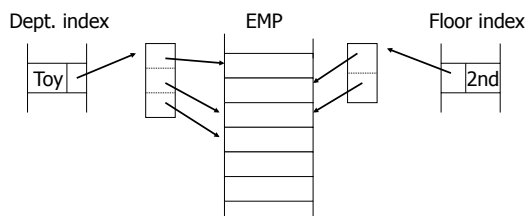
Indexes	Records
Name: primary	EMP (name,dept,floor,...)
Dept: secondary	
Floor: secondary	

6: B+trees

Data Engineering

50

Query: Get employees in  
(Toy Dept)  $\wedge$  (2nd floor)



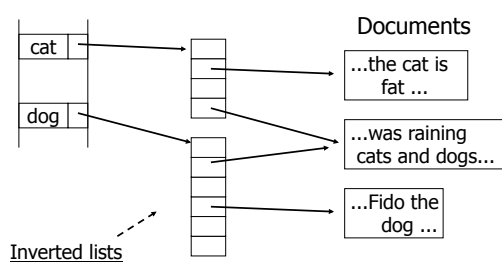
→ Intersect toy bucket and 2nd Floor bucket to get set of matching EMP's

6: B+trees

Data Engineering

51

This idea used in  
text information retrieval



6: B+trees

Data Engineering

52