# Introduction to Database Systems
## CSE 444

### Lectures 19:
### Data Storage and Indexes

November 14, 2007

1

# Outline

- Representing data elements (12)
- Index structures (13.1, 13.2)
- B-trees (13.3)

2

# Files and Tables

- A disk = a sequence of blocks

- A file = a subsequence of blocks, usually contiguous

- Need to store tables/records/indexes in files/block

3

# Representing Data Elements

- Relational database elements:

  CREATE TABLE Product (

      pid INT PRIMARY KEY,
      name CHAR(20),
      description VARCHAR(200),
      maker CHAR(10) REFERENCES Company(name)
  )

- A tuple is represented as a record
- The table is a sequence of records
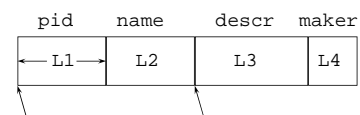
4

# Issues

- Represent attributes inside the records
- Represent the records inside the blocks

5

# Record Formats: Fixed Length

| pid | name | descr | maker |
|-----|------|-------|-------|
| L1  | L2   | L3    | L4    |

Base address (B)    Address = B+L1+L2

- Information about field types same for all records in a file; stored in *system catalogs.*
- Finding $i'th$ field requires scan of record.
- **Note the importance of schema information!**

6

## Record Header

To schema

length

| pid | name | descr | maker |
|-----|------|-------|-------|
| ←L1→ | L2 | L3 | L4 |

header

timestamp

Need the header because:
•The schema may change
   for a while new+old may coexist
•Records from different relations may coexist

7

## Variable Length Records

Other header information

header

| pid | name | descr | maker |
|-----|------|-------|-------|
| ←L1→ | L2 | L3 | L4 |

length

Place the fixed fields first: F1
Then the variable length fields: F2, F3, F4
Null values take 2 bytes only
Sometimes they take 0 bytes (when at the end)

8

## Storing Records in Blocks

• Blocks have fixed size (typically 4k – 8k)

BLOCK

| | | R4 | R3 | R2 | R1 |
|---|---|----|----|----|----|

9

## BLOB

• Binary large objects
• Supported by modern database systems
• E.g. images, sounds, etc.
• Storage: attempt to cluster blocks together

CLOB = character large object
• Supports only restricted operations

10
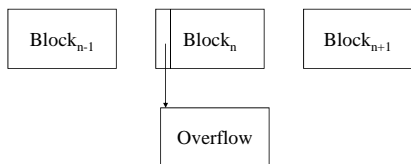
## File Types

• Unsorted (heap)

• Sorted (e.g. by pid)

11

## Modifications: Insertion

• File is unsorted: add it to the end (easy ☺)
• File is sorted:
  – Is there space in the right block ?
    • Yes: we are lucky, store it there
  – Is there space in a neighboring block ?
    • Look 1-2 blocks to the left/right, shift records
  – If anything else fails, create *overflow block*

12

## Overflow Blocks

| Block$_{n-1}$ | Block$_n$ | Block$_{n+1}$ |
|---|---|---|

Overflow

- After a while the file starts being dominated by overflow blocks: time to reorganize

13

## Modifications: Deletions

- Free space in block, shift records

- May be able to eliminate an overflow block

- Can never really eliminate the record, because others may *point* to it
  - Place a tombstone instead (a NULL record)

How can we *point* to a record in an RDBMS ?    14

## Modifications: Updates

- If new record is shorter than previous, easy ☺
- If it is longer, need to shift records, create overflow blocks

15

## Pointers

Logical pointer to a record consists of:

- Logical block number

- An offset in the block's header

We use pointers in Indexes and in Log entries

Note: review what a pointer in C is    16

## Indexes

- An *index* on a file speeds up selections on the *search key fields* for the index.
  - Any subset of the fields of a relation can be the search key for an index on the relation.
  - *Search key* is not the same as *key* (minimal set of fields that uniquely identify a record in a relation).
- An index contains a collection of *data entries*, and supports efficient retrieval of all data entries with a given key value **k**.
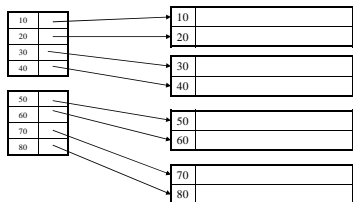
## Index Classification

- Clustered/unclustered
  - Clustered = records close in the index are close in the data; same as saying that the table is ordered by the index key
  - Unclustered = records close in the index may be far in the data
- Primary/secondary
  - Interpretation 1:
    - Primary = is over attributes part of the primary
    - Secondary = cannot reorder data
  - Interpretation 2: means the same as clustered/unclustured
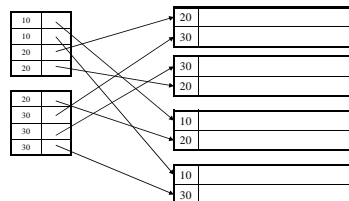- B+ tree or Hash table

18

## Clustered Index

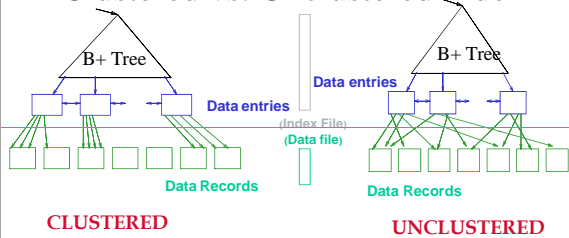- File is sorted on the index attribute
- Only one per table



19

## Unclustered Index

- Several per table



20

## Clustered vs. Unclustered Index
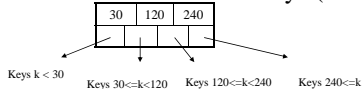


**CLUSTERED**          **UNCLUSTERED**

## B+ Trees

- Search trees
- Idea in B   Trees:
  - make 1 node = 1 block
- Idea in B+ Trees:
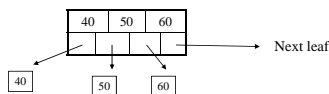  - Make leaves into a linked list (range queries are easier)

22

## B+ Trees Basics

- Parameter d = the *degree*
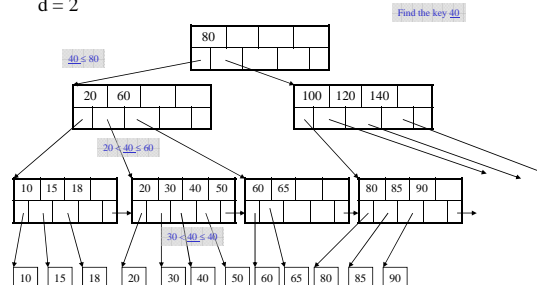- Each node has >= d and <= 2d keys (except root)



Keys k < 30      Keys 30<=k<120    Keys 120<=k<240    Keys 240<=k

- Each leaf has >=d and <= 2d keys:



Next leaf

23

## B+ Tree Example

d = 2

Find the key 40



24

## B+ Tree Design

- How large d ?
- Example:
  - Key size = 4 bytes
  - Pointer size = 8 bytes
  - Block size = 4096 byes
- 2d x 4 + (2d+1) x 8 <= 4096
- d = 170

25

## Searching a B+ Tree

- Exact key values:
  - Start at the root
  - Proceed down, to the leaf

  Select name
  From people
  Where age = 25

- Range queries:
  - As above
  - Then sequential traversal

  Select name
  From people
  Where 20 <= age
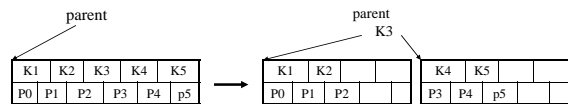   and age <= 30

26

## B+ Trees in Practice

- Typical order: 100. Typical fill-factor: 67%.
  - average fanout = 133
- Typical capacities:
  - Height 4: $133^4$ = 312,900,700 records
  - Height 3: $133^3$ =    2,352,637 records
- 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

## Insertion in a B+ Tree

Insert (K, P)
- Find leaf where K belongs, insert
- If no overflow (2d keys or less), halt
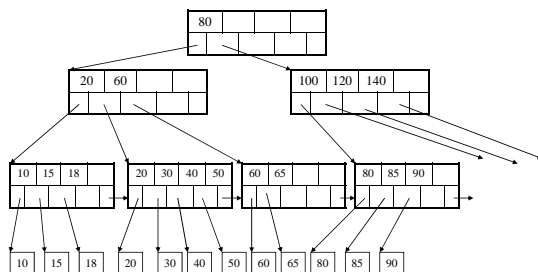- If overflow (2d+1 keys), split node, insert in parent:



- If leaf, keep K3 too in right node
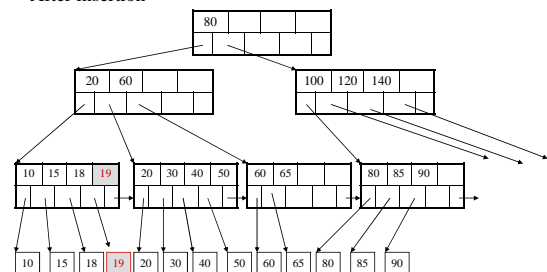- When root splits, new root has 1 key only
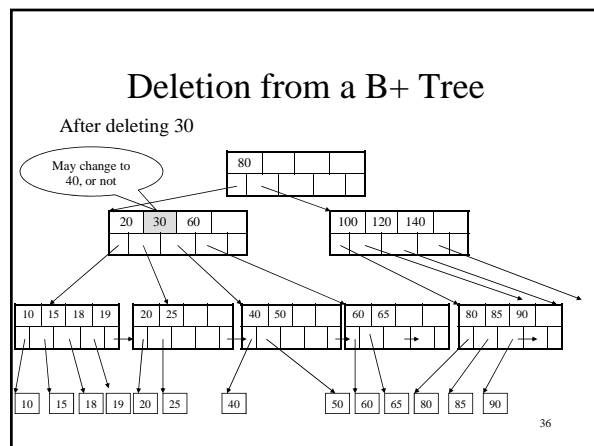
28

## Insertion in a B+ Tree

Insert K=19



29

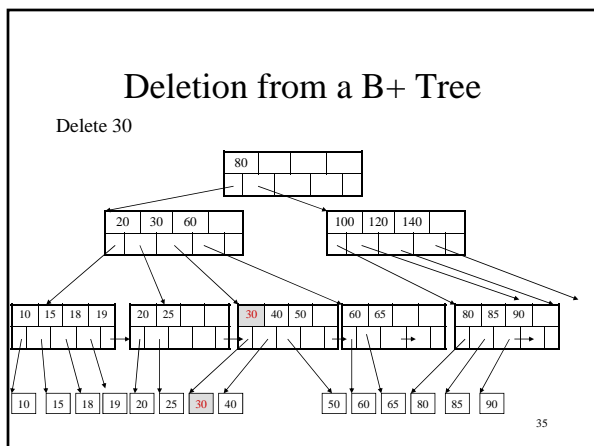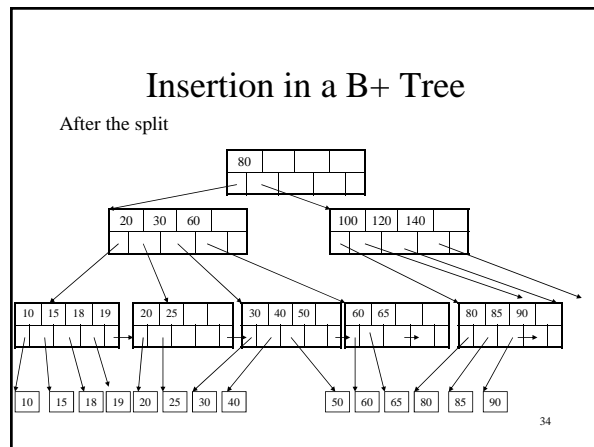## Insertion in a B+ Tree

After insertion



30

5

## Insertion in a B+ Tree

Now insert 25

80

20  60    100  120  140

10 15 18 19   20 30 40 50   60 65   80 85 90

10  15  18  19  20  30  40  50  60  65  80  85  90

31

## Insertion in a B+ Tree

After insertion

80

20  60    100  120  140

10 15 18 19   20 25 30 40 50   60 65   80 85 90

10  15  18  19  20  25  30  40  50  60  65  80  85  90

32

## Insertion in a B+ Tree

But now have to split !

80

20  60    100  120  140

10 15 18 19   20 25 30 40 50   60 65   80 85 90

10  15  18  19  20  25  30  40  50  60  65  80  85  90

33

## Insertion in a B+ Tree

After the split

80

20  30  60    100  120  140

10 15 18 19   20 25   30 40 50   60 65   80 85 90

10  15  18  19  20  25  30  40   50  60  65  80  85  90

34

## Deletion from a B+ Tree
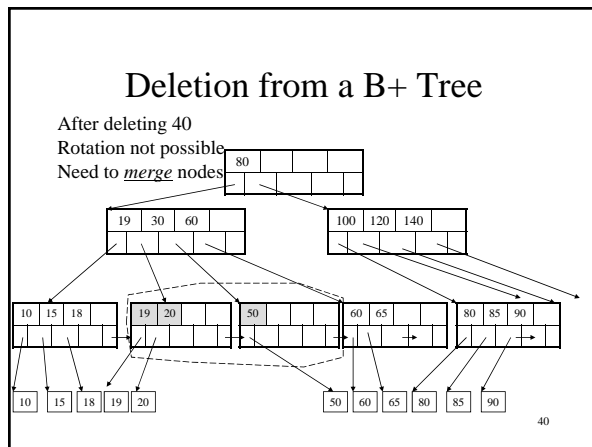
Delete 30

80

20  30  60    100  120  140

10 15 18 19   20 25   30 40 50   60 65   80 85 90

10  15  18  19  20  25  30  40   50  60  65  80  85  90

35

## Deletion from a B+ Tree

After deleting 30

May change to
40, or not

80

20  30  60    100  120  140

10 15 18 19   20 25   40 50   60 65   80 85 90

10  15  18  19  20  25   40   50  60  65  80  85  90

36

## Deletion from a B+ Tree

Now delete 25

37

## Deletion from a B+ Tree

After deleting 25
Need to rebalance
*Rotate*

38

## Deletion from a B+ Tree

Now delete 40

39

## Deletion from a B+ Tree

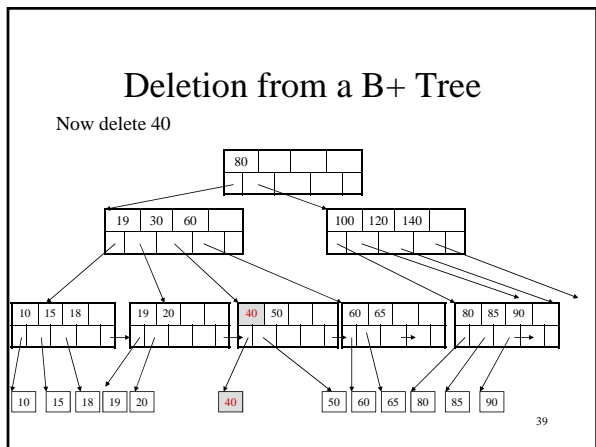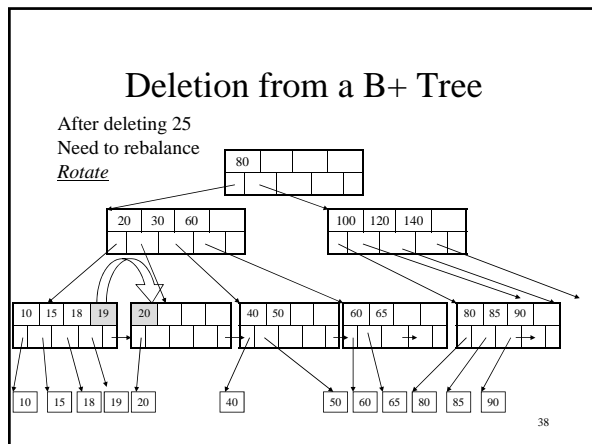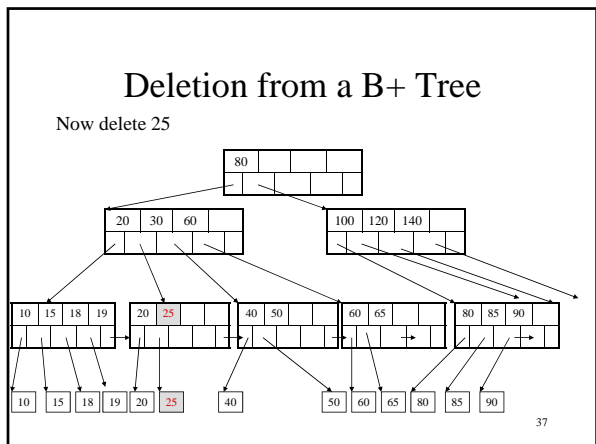After deleting 40
Rotation not possible
Need to *merge* nodes

40

## Deletion from a B+ Tree

Final tree

41

## Summary on B+ Trees

- Default index structure on most DBMS
- Very effective at answering 'point' queries:
    productName − 'gizmo'
- Effective for range queries:
    50 < price AND price < 100
- Less effective for multirange:
    50 < price < 100  AND 2 < quant < 20

42