## Slide 1

Database System Internals

# Indexing

Paul G. Allen School of Computer Science and Engineering
University of Washington, Seattle

January 15, 2021 — CSE 444 - Winter 2020 — 1

1

## Slide 2

# Announcements

- Homework 1:
  - Due on gradescope 11pm tonight

- Lab 1 complete:
  - Due on Jan. 20th

- 544 reading 1:
  - Due approx. this weekend (these due dates are flexible) by email to me

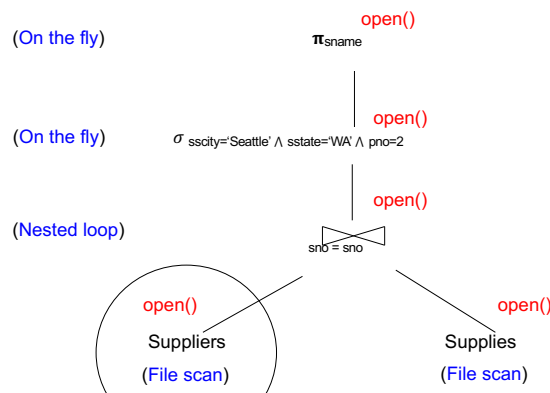January 15, 2021 — CSE 444 - Winter 2020 — 2

2

## Slide 3

# Heap File Access Method API

- **Create** or **destroy** a file
- **Insert** a record
- **Delete** a record with a given rid (rid)
  - rid: unique tuple identifier (more later)
- **Get** a record with a given rid
  - Not necessary for sequential scan operator
  - But used with indexes (more next lecture)
- **Scan** all records in the file

January 15, 2021 — CSE 444 - Winter 2020 — 3

3

## Slide 4

# Query Execution   How it all Fits

(On the fly)    $\pi_{sname}$    open()

(On the fly)    $\sigma_{sscity='Seattle' \wedge sstate='WA' \wedge pno=2}$    open()
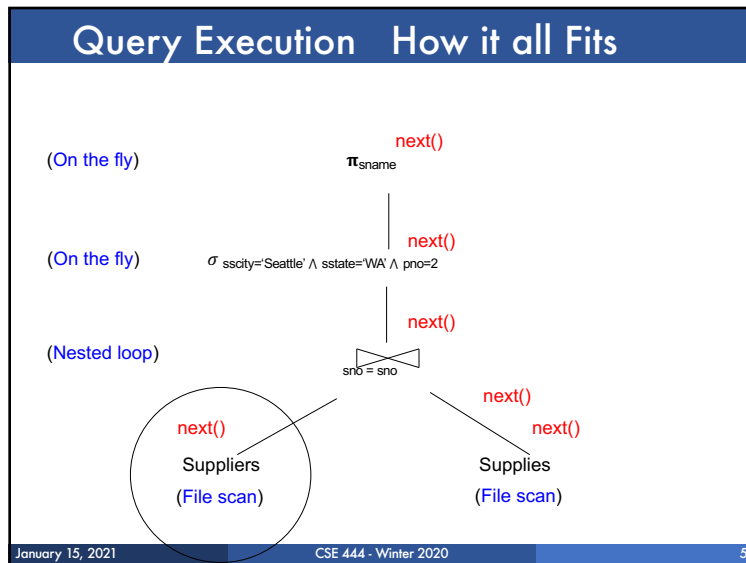
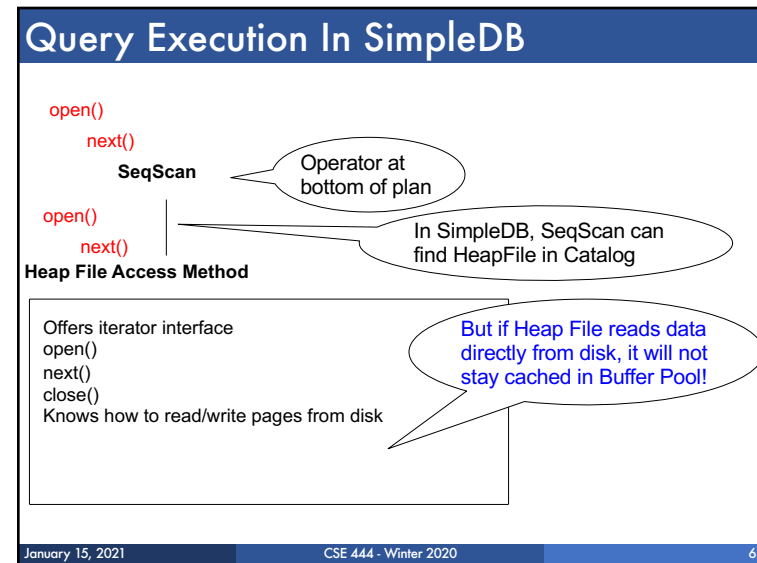(Nested loop)    sno = sno    open()

open()    Suppliers    (File scan)

open()    Supplies    (File scan)
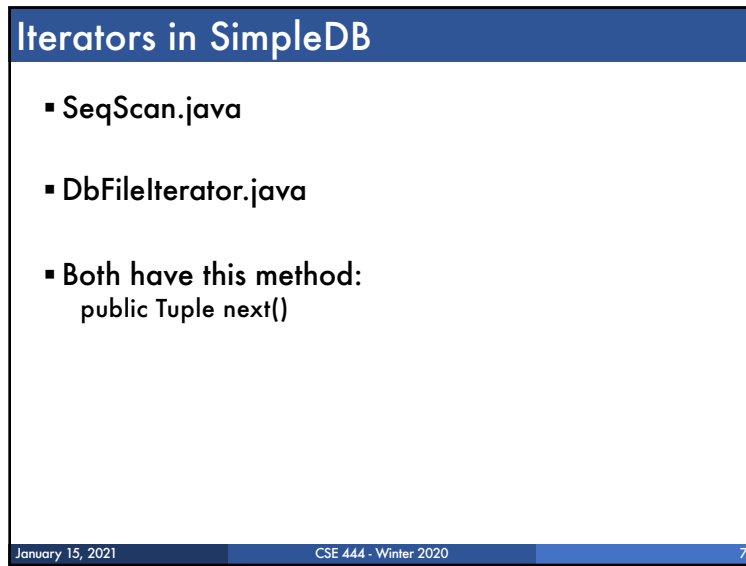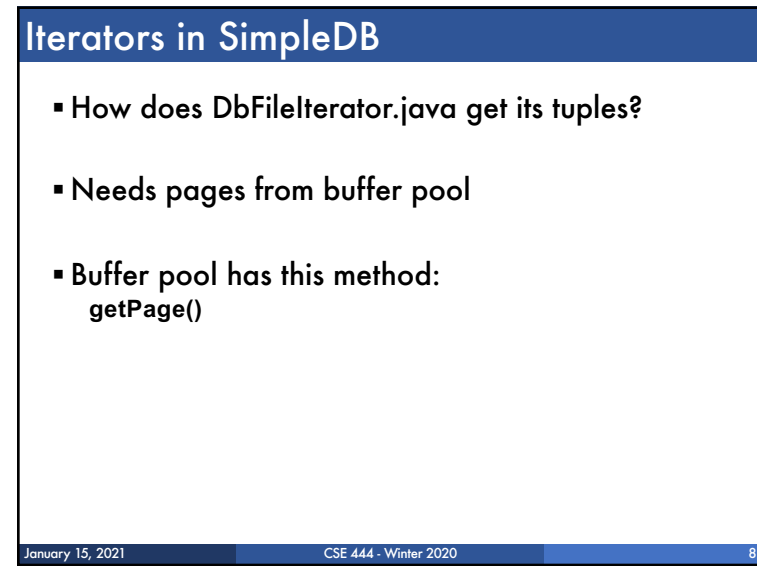
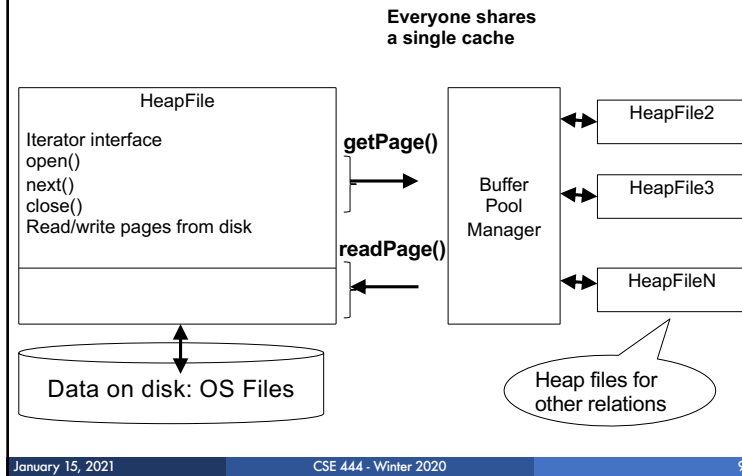January 15, 2021 — CSE 444 - Winter 2020 — 4

4

1

## Query Execution   How it all Fits

(On the fly)     next()
          π sname

(On the fly)     next()
          σ sscity='Seattle' ∧ sstate='WA' ∧ pno=2

(Nested loop)     next()
               sno = sno

          next()              next()
                       next()
     Suppliers          Supplies
     (File scan)        (File scan)

5

## Query Execution In SimpleDB

open()
     next()
          SeqScan          Operator at bottom of plan

open()
     next()          In SimpleDB, SeqScan can find HeapFile in Catalog
Heap File Access Method

Offers iterator interface          But if Heap File reads data directly from disk, it will not stay cached in Buffer Pool!
open()
next()
close()
Knows how to read/write pages from disk

6

## Iterators in SimpleDB

▪ SeqScan.java

▪ DbFileIterator.java

▪ Both have this method:
   public Tuple next()

7

## Iterators in SimpleDB

▪ How does DbFileIterator.java get its tuples?

▪ Needs pages from buffer pool

▪ Buffer pool has this method:
   getPage()

8

## Query Execution In SimpleDB

**Everyone shares a single cache**



HeapFile

Iterator interface
open()
next()
close()
Read/write pages from disk

**getPage()**

**readPage()**

Buffer Pool Manager

HeapFile2

HeapFile3

HeapFileN

Data on disk: OS Files

Heap files for other relations

9

## Buffer Manager

- Brings pages in from memory and caches them
- Eviction policies
  - Random page (ok for SimpleDB)
  - Least-recently used
  - The "clock" algorithm
- Keeps track of which **pages are dirty**
  - A dirty page has changes not reflected on disk
  - Implementation: Each page includes a dirty bit

10

## Buffer Manager



Page requests from higher-level code

Access methods
Buffer pool manager

Buffer pool

Disk page

Free frame

Main memory

Disk is a collection of blocks

Disk

1 page corresponds to 1 disk block

11

## Basic Access Method: Heap File

API
- **Create** or **destroy** a file
- **Insert** a record
- **Delete** a record with a given rid (rid)
  - rid: unique tuple identifier (more later)
- **Get** a record with a given rid
  - Not necessary for sequential scan operator
  - But used with indexes
- **Scan** all records in the file

12

3

## But Often Also Want….

- **Scan** all records in the file that match a **predicate** of the form **attribute op value**
  - Example: Find all students with GPA > 3.5

- Critical to support such requests efficiently
  - Why read all data form disk when we only need a small fraction of that data?

- This lecture and next, we will learn how

13

## Searching in a Heap File

File is not sorted on any attribute

`Student(sid: int, age: int, …)`



| 30 | 18 … |
| 70 | 21 |

1 record

| 20 | 20 |
| 40 | 19 |

1 page

| 80 | 19 |
| 60 | 18 |

| 10 | 21 |
| 50 | 22 |

14

## Heap File Search Example

- 10,000 students
- 10 student records per page
- Total number of pages: 1,000 pages
- Find student whose sid is 80
  - Must read on average 500 pages
- Find all students older than 20
  - Must read all 1,000 pages
- Can we do better?

15

## Sequential File

File sorted on an attribute, usually on primary key

`Student(sid: int, age: int, …)`

| 10 | 21 … |
| 20 | 20 |

| 30 | 18 |
| 40 | 19 |

| 50 | 22 |
| 60 | 18 |

| 70 | 21 |
| 80 | 19 |

16

## Sequential File Example

- Total number of pages: 1,000 pages
- Find student whose sid is 80
  - Could do binary search, read $\log_2(1,000) \approx 10$ pages
- Find all students older than 20
  - Must still read all 1,000 pages
- Can we do even better?

- Note: Sorted files are inefficient for inserts/deletes

January 15, 2021     CSE 444 - Winter 2020     17

17

## Creating Indexes in SQL

CREATE TABLE   V(M int,  N varchar(20),   P int);

CREATE INDEX V1 ON V(N)

CREATE INDEX V2 ON V(P, M)

select *
from V
where P=55 and M=77

select *
from V
where P=55

January 15, 2021     CSE 444 - Winter 2020     18

18

## Outline

- Index structures
- Hash-based indexes

  } Today

- B+ trees } Next time

January 15, 2021     CSE 444 - Winter 2020     19

19

## Indexes

- **Index:** data structure that organizes data records on disk to optimize selections on the *search key fields* for the index

- An index contains a collection of *data entries*, and supports efficient retrieval of all data entries with a given search key value **k**

- Indexes are also access methods!
  - So they provide the same API as we have seen for Heap Files
  - And efficiently support scans over tuples matching predicate on search key

Index File
Search key: age

| 18 |
| 18 |
| 19 |
| 19 |

| 20 |
| 21 |
| 21 |
| 22 |

| 10 | 21 |
| 20 | 20 |

| 30 | 18 |
| 40 | 19 |

| 50 | 22 |
| 60 | 18 |

| 70 | 21 |
| 80 | 19 |

Data File
(sequential file sorted on sid)

January 15, 2021     CSE 444 - Winter 2020     20

20

5

## Indexes

- **Search key** = can be any set of fields
  - not the same as the primary key, nor a key
- **Index** = collection of data entries
- **Data entry** for key k can be:
  - (k, RID)
  - (k, list-of-RIDs)
  - The actual record with key k
    - In this case, **the index is also a special file organization**
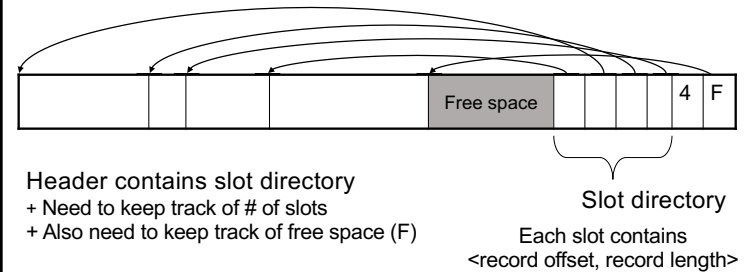    - Called: "indexed file organization"

21

## Page Format Approach 2



Header contains slot directory
+ Need to keep track of # of slots
+ Also need to keep track of free space (F)

Slot directory

Each slot contains
<record offset, record length>

Can handle variable-length records
Can move tuples inside a page without changing RIDs
RID is (PageID, SlotID) combination

22

## Different Types of Files

- For the data inside base relations:
  - Heap file (tuples stored without any order)
  - Sequential file (tuples sorted on some attribute(s))
  - Indexed file (tuples organized following an index)
- Then we can have additional index files that store (key,rid) pairs
- Index can also be a "covering index"
  - Index contains (search key + other attributes, rid)
  - Index suffices to answer some queries

23

## Primary Index

- <u>Primary</u> index determines location of indexed records
- *Dense* index: sequence of (key,rid) pairs

24

6

## Primary Index

- *Sparse* index



Can store more search keys in same number of index files

25

## Primary Index with Duplicate Keys

- Dense index:

26

## Primary Index: Back to Example

- Let's assume all pages of index fit in memory

- Find student whose sid is 80
  - Index (dense or sparse) points directly to the page
  - Only need to read 1 page from disk.
- Find all students older than 20
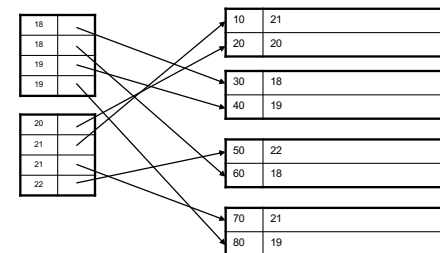
- How can we make *both* queries fast?

27

## Secondary Indexes

- Do not determine placement of records in data files
- Always dense (why ?)

28

7

## Clustered vs. Unclustered Index



**Data entries**
**(Index File)**
**(Data file)**

**Data entries**

**Data entries**

**Data Records**

**Data Records**

**CLUSTERED**

**UNCLUSTERED**

Clustered = records close in index are close in data

29

## Clustered/Unclustered

- Primary index = clustered by definition
- Secondary indexes = usually unclustered

30

## Secondary Indexes

- Applications
  - Index unsorted files (heap files)

  - When necessary to have multiple indexes

  - Index files that hold data from two relations

31

## Index Classification Summary

- Primary/secondary
  - Primary = determines the location of indexed records
  - Secondary = cannot reorder data, does not determine data location

- Dense/sparse
  - Dense = every key in the data appears in the index
  - Sparse = the index contains only some keys

- Clustered/unclustered
  - Clustered = records close in index are close in data
  - Unclustered = records close in index may be far in data
- B+ tree / Hash table / ...

32

8

## Ex1. Primary Dense Index (tid)



- **Dense:** an "index key" for every database record
  - (In this case) every "database key" appears as an "index key"
- **Primary:** determines the location of indexed records
- Also, **Clustered:** records close in index are close in data

January 15, 2021    CSE 444 - Winter 2020    33

33

## Improve from Primary Clustered Index?

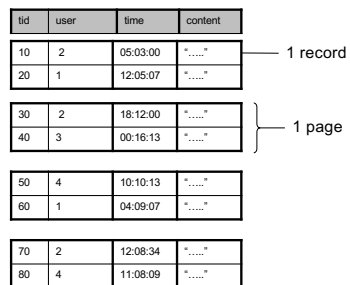Clustered Index can be made <u>Sparse</u>
(normally one key per page)

January 15, 2021    CSE 444 - Winter 2020    34

34

## Ex2. Draw a <u>primary</u> <u>sparse</u> index on "tid"



January 15, 2021    CSE 444 - Winter 2020    35

35

## Ex2. Primary Sparse Index (tid)



- Only one index file page instead of two

January 15, 2021    CSE 444 - Winter 2020    36

36

9

## Large Indexes

- What if index does not fit in memory?

- Would like to index the index itself
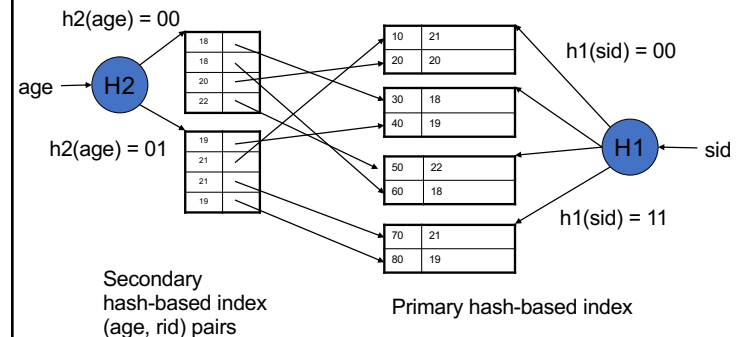  - Hash-based index
  - Tree-based index

37

## Hash-Based Index

Good for point queries but not range queries



h2(age) = 00

age → H2

h2(age) = 01

h1(sid) = 00

H1 ← sid

h1(sid) = 11

Secondary hash-based index (age, rid) pairs

Primary hash-based index

38

## Tree-Based Index

- How many index levels do we need?
- Can we create them automatically? Yes!
- Can do something even more powerful!

39

## B+ Trees

- Search trees

- Idea in B Trees
  - Make 1 node = 1 page (= 1 block)

- Idea in B+ Trees
  - Keep tree balanced in height – dynamic rather than static
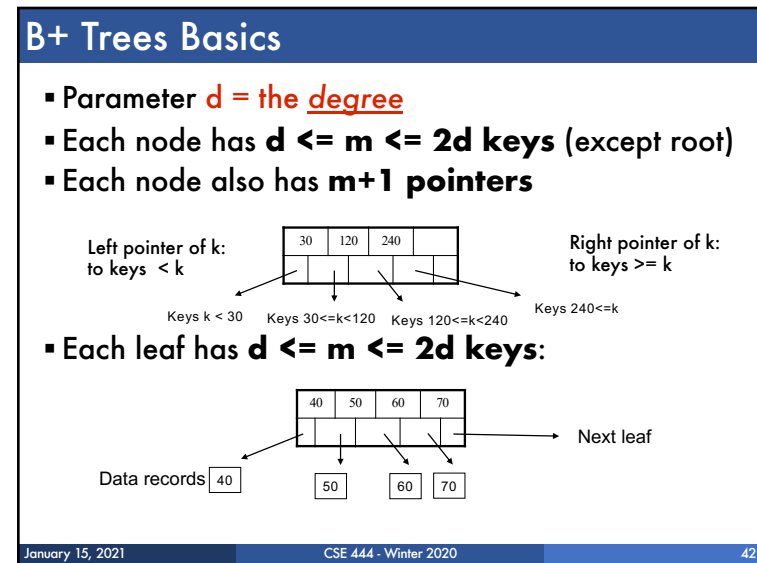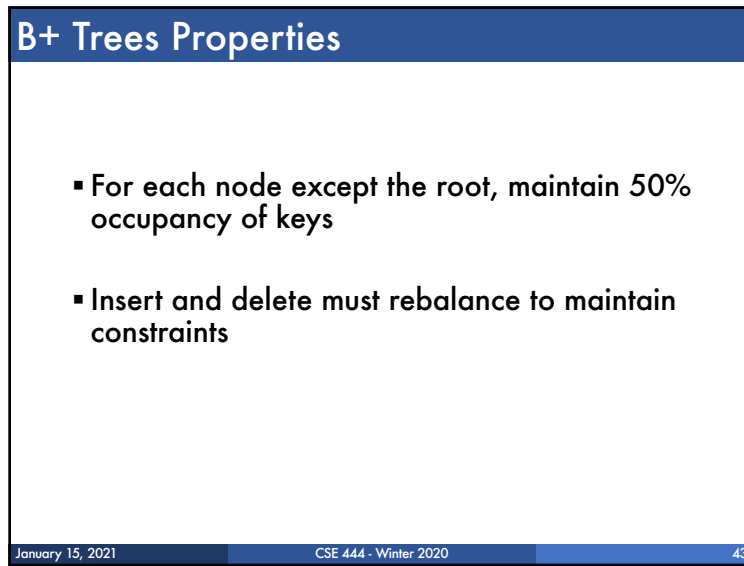  - Make leaves into a linked list : facilitates range queries

40

10

## [Slide 41]



**Data entries**

**Data entries**
**(Index File)**

**Data entries**
**(Data file)**

**Data Records**

**Data Records**

**CLUSTERED**

**UNCLUSTERED**

Note: can also store data records directly as data entries

January 15, 2021 — CSE 444 - Winter 2020 — 41

41

## B+ Trees Basics

- Parameter d = the *degree*
- Each node has **d <= m <= 2d keys** (except root)
- Each node also has **m+1 pointers**

Left pointer of k:
to keys < k

| 30 | 120 | 240 | |

Right pointer of k:
to keys >= k

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

- Each leaf has **d <= m <= 2d keys**:

| 40 | 50 | 60 | 70 |

Next leaf

Data records  40     50     60  70

January 15, 2021 — CSE 444 - Winter 2020 — 42

42

## B+ Trees Properties

- For each node except the root, maintain 50% occupancy of keys

- Insert and delete must rebalance to maintain constraints

January 15, 2021 — CSE 444 - Winter 2020 — 43

43

## Searching a B+ Tree

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

- Range queries:
  - Find lowest bound as above
  - Then sequential traversal

Select name
From Student
Where age = 25

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

January 15, 2021 — CSE 444 - Winter 2020 — 44

44

11

## B+ Tree Example

d = 2

Find the key 40

40 < 80

| 80 | | |
|---|---|---|

20 < 40 < 60

| 20 | 60 | |

| 100 | 120 | 140 |

| 10 | 15 | 18 |

| 20 | 30 | 40 | 50 |

| 60 | 65 | |

| 80 | 85 | 90 |

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

45

## B+ Tree Design

- How large d ?    Make one node fit on one block
- Example:
  - Key size = 4 bytes
  - Pointer size = 8 bytes
  - Block size = 4096 bytes

| 30 | 120 | 240 | |
|---|---|---|---|

(e.g. d = 2)

- $2d \times 4 + (2d+1) \times 8 <= 4096$
- d = 170

46

## 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

47

## 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:

parent

parent
K3

| K1 | K2 | K3 | K4 | K5 |
|---|---|---|---|---|
| P0 | P1 | P2 | P3 | P4 | P5 |

| K1 | K2 | | |
|---|---|---|---|
| P0 | P1 | P2 | |

| K4 | K5 | |
|---|---|---|
| P3 | P4 | p5 |

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

48

12

## Insertion in a B+ Tree

Insert K=19

49

## Insertion in a B+ Tree

After insertion

50

## Insertion in a B+ Tree

Now insert 25

51

## Insertion in a B+ Tree

After insertion

52

13

## Insertion in a B+ Tree

But now have to split !

53

## Insertion in a B+ Tree

After the split

54

## Deletion in a B+ Tree

Delete (K, P)

- Find leaf where K belongs, delete
- Check for capacity
- If leaf below capacity, search adjacent nodes (left first, then right) for extra tuples and rotate them to new leaf
- If adjacent nodes at 50% full, merge
- Update and repeat algorithm on parent nodes if necessary

55

## Deletion from a B+ Tree

Delete 30

56

14

## Deletion from a B+ Tree

After deleting 30



## Deletion from a B+ Tree

Now delete 25



## Deletion from a B+ Tree

After deleting 25
Need to rebalance
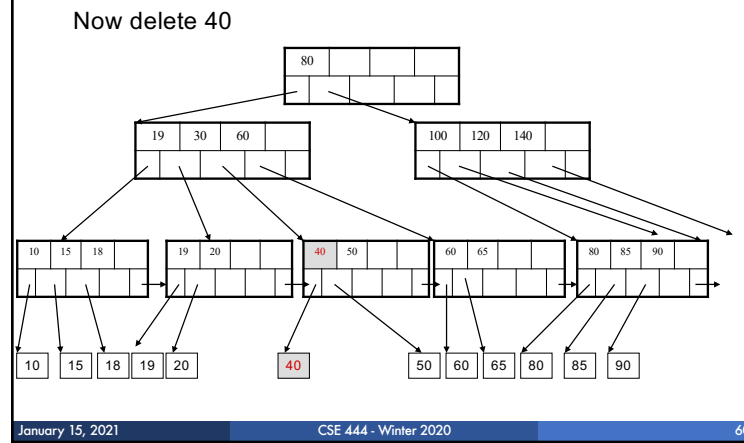*Rotate*



## Deletion from a B+ Tree

Now delete 40

## Deletion from a B+ Tree

After deleting 40
Rotation not possible
Need to _merge_ nodes

61

## Deletion from a B+ Tree

Final tree

62

## Summary on B+ Trees

- Default index structure on most DBMSs
- 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

63

16