# Tree-Structured Indexes

### Chapter 10

1

---

## Alternative File Organizations

Many alternatives exist, *each ideal for some situations, and not so good in others:*

- Heap (random order) files:  Suitable when typical access is a file scan retrieving all records.
- Sorted Files:  Best if records must be retrieved in some order, or only a `range' of records is needed.
- Indexes: Data structures to organize records via trees or hashing.
  - Like sorted files, they speed up searches for a subset of records, based on values in certain ("search key") fields
  - Updates are much faster than in sorted files.

2

---

## 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 **k\*** with a given key value **k**.
  - Given data entry k\*, we can find record with key k in at most one disk I/O.  (Details soon …)

3

### How the index relates to the underlying table

❖ *As for any index, 3 alternatives for data entries* **k\***:
  ▪ Data record with key value **k**
  ▪ <**k**, rid of data record with search key value **k**>
  ▪ <**k**, list of rids of data records with search key **k**>

❖ Choice is orthogonal to the *indexing technique* used to locate data entries **k\***.

❖ Tree-structured indexing techniques support both *range searches* and *equality searches*.

❖ <u>*ISAM*</u>: static structure; <u>*B+ tree*</u>: dynamic, adjusts gracefully under inserts and deletes.

4

---

### Alternatives for Data Entries (Contd.)

❖ Alternative 1:
  ▪ If this is used, index structure is a file organization for data records (instead of a Heap file or sorted file).
  ▪ At most one index on a given collection of data records can use Alternative 1.  (Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency.)
  ▪ If data records are very large,  # of pages containing data entries is high.  Implies size of auxiliary information in the index is also large, typically.

5

---

### Alternatives for Data Entries (Contd.)

❖ Alternatives 2 and 3:
  ▪ Data entries typically much smaller than data records.  So, better than Alternative 1 with large data records, especially if search keys are small. (Portion of index structure used to direct search, which depends on size of data entries, is much smaller than with Alternative 1.)
  ▪ Alternative 3 more compact than Alternative 2, but leads to variable sized data entries even if search keys are of fixed length.

6

## *Index Classification*

- ❖ *Primary* vs. *secondary*: If search key contains primary key, then called primary index.
  - *Unique* index: Search key contains a candidate key.
- ❖ *Clustered* vs. *unclustered*: If order of data records is the same as, or `close to', order of data entries, then called clustered index.
  - Alternative 1 implies clustered; in practice, clustered also implies Alternative 1 (since sorted files are rare).
  - A file can be clustered on at most one search key.
  - Cost of retrieving data records through index varies *greatly* based on whether index is clustered or not!
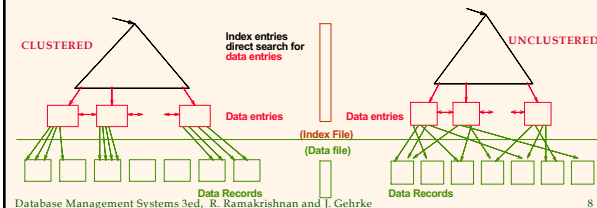
7

## *Clustered vs. Unclustered Index*

- ❖ Suppose that Alternative (2) is used for data entries, and that the data records are stored in a Heap file.
  - To build clustered index, first sort the Heap file (with some free space on each page for future inserts).
  - Overflow pages may be needed for inserts. (Thus, order of data recs is `close to', but not identical to, the sort order.)



**CLUSTERED**   **Index entries direct search for data entries**   **UNCLUSTERED**

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

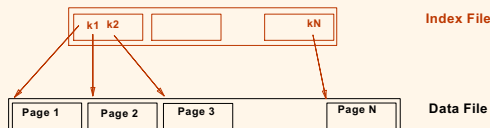**Data Records**   **Data Records**

8

## *Range Searches*

- ❖ ``*Find all students with gpa > 3.0*''
  - If data is in sorted file, do binary search to find first such student, then scan to find others.
  - Cost of binary search can be quite high.
- ❖ Simple idea: Create an `index' file.



| k1 | k2 | | | kN | **Index File** |

| Page 1 | Page 2 | Page 3 | | Page N | **Data File** |

☞ *Can do binary search on (smaller) index file!*

9

## ISAM

**index entry**

| $P_0$ | $K_1$ | $P_1$ | $K_2$ | $P_2$ | $\diamond$ $\diamond$ $\diamond$ | $K_m$ | $P_m$ |

❖ Index file may still be quite large. But we can apply the idea repeatedly!

**Non-leaf Pages**

- - -

**Leaf Pages**

- - -        - - -

**Overflow page**

**Primary pages**

☛ *Leaf pages contain data entries.*

10

---

## Comments on ISAM

| Data Pages |
|---|
| **Index Pages** |
| **Overflow pages** |

❖ *File creation*: Leaf (data) pages allocated sequentially, sorted by search key; then index pages allocated, then space for overflow pages.

❖ *Index entries*: <search key value, page id>; they `direct' search for *data entries*, which are in leaf pages.

❖ *Search*: Start at root; use key comparisons to go to leaf. Cost $\propto \log_F N$ ; F = # entries/index pg, N = # leaf pgs

❖ *Insert*: Find leaf data entry belongs to, and put it there.

❖ *Delete*: Find and remove from leaf; if empty overflow page, de-allocate.

☛ **Static tree structure**: *inserts/deletes affect only leaf pages*.

11

---

## Example ISAM Tree

❖ Each node can hold 2 entries; no need for `next-leaf-page' pointers. (Why?)

**Root**

| 40 | |

| 20 | 33 |          | 51 | 63 |

| 10* | 15* | | 20* | 27* | | 33* | 37* | | 40* | 46* | | 51* | 55* | | 63* | 97* |

12

## After Inserting 23*, 48*, 41*, 42*...

**Index Pages**

Root

40

20 33

51 63

**Primary Leaf Pages**

10* 15* | 20* 27* | 33* 37* | 40* 46* | 51* 55* | 63* 97*

**Overflow Pages**

23*

48* 41*

42*

13

---

## ... Then Deleting 42*, 51*, 97*

Root

40

20 33

51 63

10* 15* | 20* 27* | 33* 37* | 40* 46* | 55* | 63*

23*

48* 41*

☛ *Note that 51\* appears in index levels, but not in leaf!*

14

---

## B+ Tree: Most Widely Used Index

❖ Insert/delete at $\log_F N$ cost; keep tree *height-balanced*. (F = fanout, N = # leaf pages)

❖ Minimum 50% occupancy (except for root). Each node contains $d <= m <= 2d$ entries. The parameter **d** is called the *order* of the tree.

❖ Supports equality and range-searches efficiently.

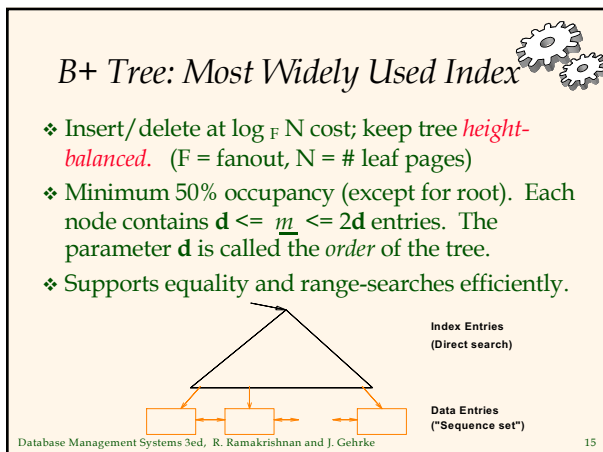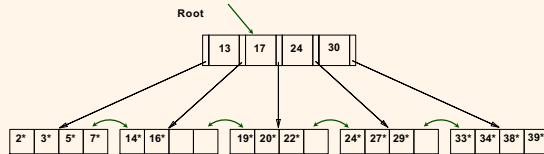**Index Entries**
(Direct search)

**Data Entries**
("Sequence set")

15

## Example B+ Tree

- ❖ Search begins at root, and key comparisons direct it to a leaf (as in ISAM).
- ❖ Search for 5*, 15*, all data entries >= 24* ...

**Root**

| 13 | 17 | 24 | 30 |

| 2* | 3* | 5* | 7* | | 14* | 16* | | | | 19* | 20* | 22* | | 24* | 27* | 29* | | 33* | 34* | 38* | 39* |

☛ *Based on the search for 15\*, we know it is not in the tree!*

16

In this example, 13,17,30 are deleted from leaf level

---

## 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

17

---

## Inserting a Data Entry into a B+ Tree

- ❖ Find correct leaf *L.*
- ❖ Put data entry onto *L.*
  - If *L* has enough space, *done*!
  - Else, must *split* *L (into L and a new node L2)*
    - Redistribute entries evenly, **copy up** middle key.
    - Insert index entry pointing to *L2* into parent of *L.*
- ❖ This can happen recursively
  - To split index node, redistribute entries evenly, but **push up** middle key.  (Contrast with leaf splits.)
- ❖ Splits "grow" tree; root split increases height.
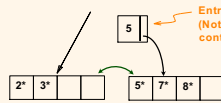  - Tree growth: gets *wider* or *one level taller at top.*

18

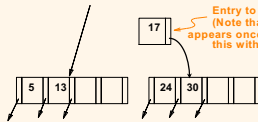## Inserting 8* into Example B+ Tree

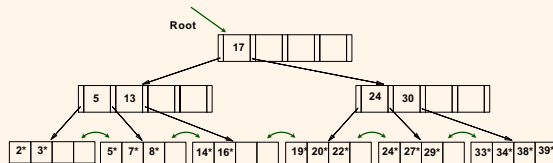- ❖ Observe how minimum occupancy is guaranteed in both leaf and index pg splits.
- ❖ Note difference between *copy-up* and *push-up*; be sure you understand the reasons for this.

Entry to be inserted in parent node. (Note that 5 is copied up and continues to appear in the leaf.)

```
      5
2*  3*      5*  7*  8*
```

Entry to be inserted in parent node. (Note that 17 is pushed up and only appears once in the index. Contrast this with a leaf split.)

```
       17
5  13        24  30
```

---

## Example B+ Tree After Inserting 8*

Root

```
              17
    5  13              24  30
2* 3*   5* 7* 8*  14* 16*   19* 20* 22*   24* 27* 29*   33* 34* 38* 39*
```

- ❖ Notice that root was split, leading to increase in height.
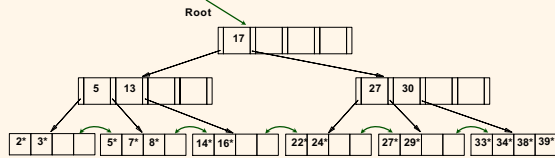- ❖ In this example, we can avoid split by re-distributing entries; however, this is usually not done in practice.

---

## Deleting a Data Entry from a B+ Tree

- ❖ Start at root, find leaf *L* where entry belongs.
- ❖ Remove the entry.
  - If L is at least half-full, *done!*
  - If L has only **d-1** entries,
    - Try to re-distribute, borrowing from *sibling (adjacent node with same parent as L).*
    - If re-distribution fails, *merge L* and sibling.
- ❖ If merge occurred, must delete entry (pointing to *L* or sibling) from parent of *L*.
- ❖ Merge could propagate to root, decreasing height.

## Example Tree After (Inserting 8*, Then) Deleting 19* and 20* ...

Root

| 17 | | | |

| 5 | 13 | | |    | 27 | 30 | | |

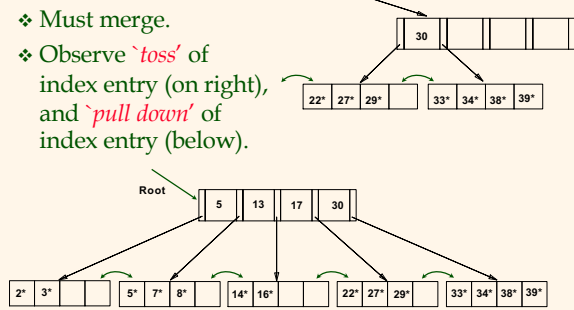| 2* | 3* | |    | 5* | 7* | 8* |    | 14* | 16* |    | 22* | 24* | |    | 27* | 29* | |    | 33* | 34* | 38* | 39* |

❖ Deleting 19* is easy.
❖ Deleting 20* is done with re-distribution. Notice how middle key is *copied up*.

22

If merging will cause overflow, we borrow an element from sibling, and update parent

---

## ... And Then Deleting 24*

❖ Must merge.
❖ Observe `*toss*' of index entry (on right), and `*pull down*' of index entry (below).

| 30 | | | |

| 22* | 27* | 29* |    | 33* | 34* | 38* | 39* |

Root

| 5 | 13 | 17 | 30 |

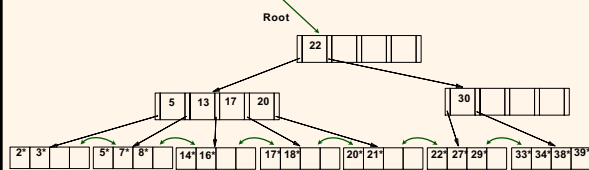| 2* | 3* | |    | 5* | 7* | 8* |    | 14* | 16* | |    | 22* | 27* | 29* |    | 33* | 34* | 38* | 39* |

23

In case if your sibling has less than 50% occupancy if being borrowed, we can merge two nodes.
Since the right side only has 30 (less than 50% occupancy) and the left side is on the boundary of 50% occupancy, we will merge again, and the result looks like the bottom image
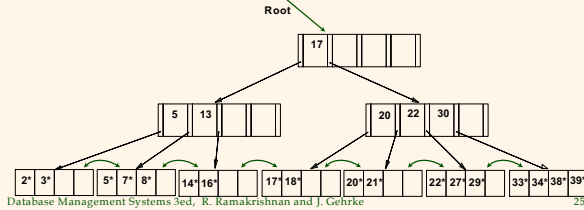
---

## Example of Non-leaf Re-distribution

❖ Tree is shown below *during deletion* of 24*. (What could be a possible initial tree?)
❖ In contrast to previous example, can re-distribute entry from left child of root to right child.

Root

| 22 | | | |

| 5 | 13 | 17 | 20 |    | 30 | | |

| 2* | 3* | |    | 5* | 7* | 8* |    | 14* | 16* |    | 17* | 18* | |    | 20* | 21* | |    | 22* | 27* | 29* |    | 33* | 34* | 38* | 39* |

24

## After Re-distribution

- Intuitively, entries are re-distributed by `pushing through` the splitting entry in the parent node.
- It suffices to re-distribute index entry with key 20; we've re-distributed 17 as well for illustration.

Root
```
        17
   5  13          20  22  30

2* 3*   5* 7* 8*  14*16*   17*18*  20*21*  22*27*29*  33*34*38*39*
```

25

just for exercise, moving one element will be enough

B+ tree is non-trival because you have to think about concurrency and recovery, etc.

## Prefix Key Compression

- Important to increase fan-out. (Why?)
- Key values in index entries only `direct traffic`; can often compress them.
  - E.g., If we have adjacent index entries with search key values *Dannon Yogurt*, *David Smith* and *Devarakonda Murthy*, we can abbreviate *David Smith* to *Dav*. (The other keys can be compressed too ...)
    - Is this correct? Not quite! What if there is a data entry *Davey Jones*? (Can only compress *David Smith* to *Davi*)
    - In general, while compressing, must leave each index entry greater than every key value (in any subtree) to its left.
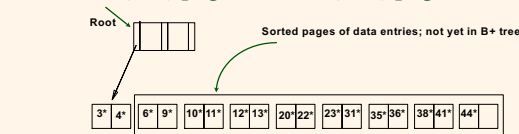- Insert/delete must be suitably modified.

26

You can take the prefix to distinguish items. For example, first three letters are enough to distinguish the three names in e.g.

## how to do more than 1 item a time?
## Bulk Loading of a B+ Tree

- If we have a large collection of records, and we want to create a B+ tree on some field, doing so by repeatedly inserting records is very slow.
- *Bulk Loading* can be done much more efficiently.
- *Initialization*: Sort all data entries, insert pointer to first (leaf) page in a new (root) page.

Root
```
Sorted pages of data entries; not yet in B+ tree

3* 4*  6* 9*  10*11*  12*13*  20*22*  23*31*  35*36*  38*41*  44*
```

27

scenario: you have a chunk of data and you want to load them all to B+ tree at a time.
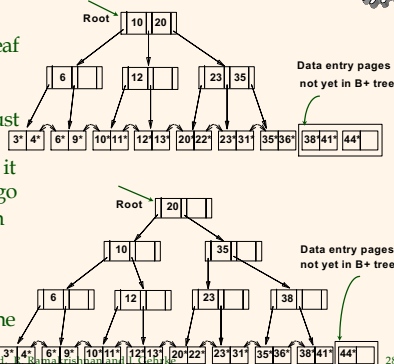
## Bulk Loading (Contd.)

- ❖ Index entries for leaf pages always entered into right-most index page just above leaf level. When this fills up, it splits. (Split may go up right-most path to the root.)
- ❖ Much faster than repeated inserts, especially when one considers locking!

Root  10  20

Data entry pages not yet in B+ tree

6    12    23  35

3* 4*  6* 9*  10*11*  12*13*  20*22*  23*31*  35*36*  38*41*  44*

Root  20

Data entry pages not yet in B+ tree

10    35

6    12    23    38

3* 4*  6* 9*  10*11*  12*13*  20*22*  23*31*  35*36*  38*41*  44*

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke    28

28

The ~~complexity is still nlogn because~~ its pre-~~sorted~~

---

## Summary of Bulk Loading

- ❖ Option 1: multiple inserts.
  - ▪ Slow.
  - ▪ Does not give sequential storage of leaves.
- ❖ Option 2: *Bulk Loading*
  - ▪ Has advantages for concurrency control.
  - ▪ Fewer I/Os during build.
  - ▪ Leaves will be stored sequentially (and linked, of course).
  - ▪ Can control "fill factor" on pages.

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke    29

29

## Composite Search Keys

- ❖ To retrieve Emp records with *age*=30 AND *sal*=4000, an index on <*age,sal*> would be better than an index on *age* or an index on *sal*.
  - ▪ Choice of index key orthogonal to clustering etc.
- ❖ If condition is: 20<*age*<30 AND 3000<*sal*<5000:
  - ▪ Clustered tree index on <*age,sal*> or <*sal,age*> is best.
- ❖ If condition is: *age*=30 AND 3000<*sal*<5000:
  - ▪ Clustered <*age,sal*> index much better than <*sal,age*> index!
- ❖ Composite indexes are larger, updated more often.

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke    30

30

search ~~key:~~
predicate selectivity:
The highly selective predicate (that ~~eliminate most tuples) first~~
to reduce the number of I/Os

~~Intersection of tuple ids:~~
get the set of tuple ids from range in age and get the set of tuple ids from range in salary. ~~Intersect them and you get~~ the qualifying tuples (We are not doing any I/O)
a composite search key of age, sal is helpful to look for age, sal and to look for age, ~~but not sal only~~
a composite search key of sal, age is helpful to look for sal, age and to look for sal, ~~but not age only~~

canned query: query that is well optimized at ~~design time (the programmer pre-know what kind of query to send and optimize accordingly)~~
Ad hoc query: unpredictable query on the fly

plan to answer a query without touching underlying data table

## Index-Only Plans

❖ A number of queries can be answered without retrieving any tuples from one or more of the relations involved if a suitable index is available.

<E.dno>

SELECT E.dno, COUNT(*)
FROM Emp E
GROUP BY E.dno

<E.dno,E.sal>
*Tree index!*

SELECT E.dno, MIN(E.sal)
FROM Emp E
GROUP BY E.dno

<E. age,E.sal>
or
<E.sal, E.age>
*Tree index!*

SELECT AVG(E.sal)
FROM Emp E
WHERE E.age=25 AND
E.sal BETWEEN 3000 AND 5000

31

go into the leaf level, since its ordered by salary, get the first node under each dno.

The ~~more index your built, the slower insertion, deletion will be, because~~ you need to update indexes.
so we need to decide whether its search heavy or update heavy

## Index-Only Plans (Contd.)

❖ Index-only plans are possible if the key is <dno,age> or we have a tree index with key <age,dno>
  ▪ Which is better?
  ▪ What if we consider the second query?

SELECT E.dno, COUNT (*)
FROM Emp E
WHERE E.age=30
GROUP BY E.dno

SELECT E.dno, COUNT (*)
FROM Emp E
WHERE E.age>30
GROUP BY E.dno

32

## Index-Only Plans (Contd.)

❖ Index-only plans can also be found for queries involving more than one table; more on this later.

<E.dno>

SELECT D.mgr
FROM Dept D, Emp E
WHERE D.dno=E.dno

<E.dno,E.eid>

SELECT D.mgr, E.eid
FROM Dept D, Emp E
WHERE D.dno=E.dno

33

take dno in Employee index to search department table
find ~~managers whose department has employees~~

## A Note on `Order'

❖ *Order* (**d**) concept replaced by physical space criterion in practice (`*at least half-full'*).
  ▪ Index pages can typically hold many more entries than leaf pages.
  ▪ Variable sized records and search keys mean differnt nodes will contain different numbers of entries.
  ▪ Even with fixed length fields, multiple records with the same search key value (*duplicates*) can lead to variable-sized data entries (if we use Alternative (3)).

34

## Summary

❖ Tree-structured indexes are ideal for range-searches, also good for equality searches.
❖ ISAM is a static structure.
  ▪ Only leaf pages modified; overflow pages needed.
  ▪ Overflow chains can degrade performance unless size of data set and data distribution stay constant.
❖ B+ tree is a dynamic structure.
  ▪ Inserts/deletes leave tree height-balanced; $\log_F N$ cost.
  ▪ High fanout (**F**) means depth rarely more than 3 or 4.
  ▪ Almost always better than maintaining a sorted file.

35

## Summary (Contd.)

  ▪ Typically, 67% occupancy on average.
  ▪ Usually preferable to ISAM, modulo *locking* considerations; adjusts to growth gracefully.
  ▪ If data entries are data records, splits can change rids!
❖ Key compression increases fanout, reduces height.
❖ Bulk loading can be much faster than repeated inserts for creating a B+ tree on a large data set.
❖ Most widely used index in database management systems because of its versatility.  One of the most optimized components of a DBMS.

36