

# Data storage

Tuesday, May 7, 2019 6:15 PM

- Hard Disk mechanical devices

- block  $\rightarrow$  4 ~ 16 GB unit of read/write  $\Rightarrow$  indispensable; always read one block.

- Once in memory, we call it a page.

- rotation speed: not increase.

- density: slowly increased

- CPU once read 1/4/8 byte.

- Read at rotation speed

- fast: Sequential read

- slow: random read

- Rule of thumb

- random read 1% ~ 2% of file = sequential read the whole file.

- 1~2% decrease over time because of increased memory.

- Data Storage

- DBMS store data in files

- On disk, a file is split into blocks

- o  $\begin{cases} \text{Heap file} & \text{unsorted} \\ \text{Sequential file} & \text{sorted according to key} \end{cases}$

- **Index**: neither Heapfile nor sequentialfile.

- o **def**: an additional file that allows fast access to records in the data file given a search key.

- o Contains: (key, value) pairs

- key = an attribute value  $\rightarrow$  here means search key
- value: a pointer to the record or the record itself.

- Key  $\begin{cases} \text{Primary key} & \text{unique} \\ \text{key of a sequential file} & \text{how file is sorted.} \\ \text{Index key} & \text{how index is organized} \end{cases}$

- o Index can be  $\begin{cases} \text{Dense} & \text{one entry per record} \\ \text{Sparse} & \text{one entry per block} \end{cases}$

Binary search on index

e.g. 2. Index on fName

Index can only be dense. ( $\because$  in data file, data is not put in block by fName)

Index -> index file

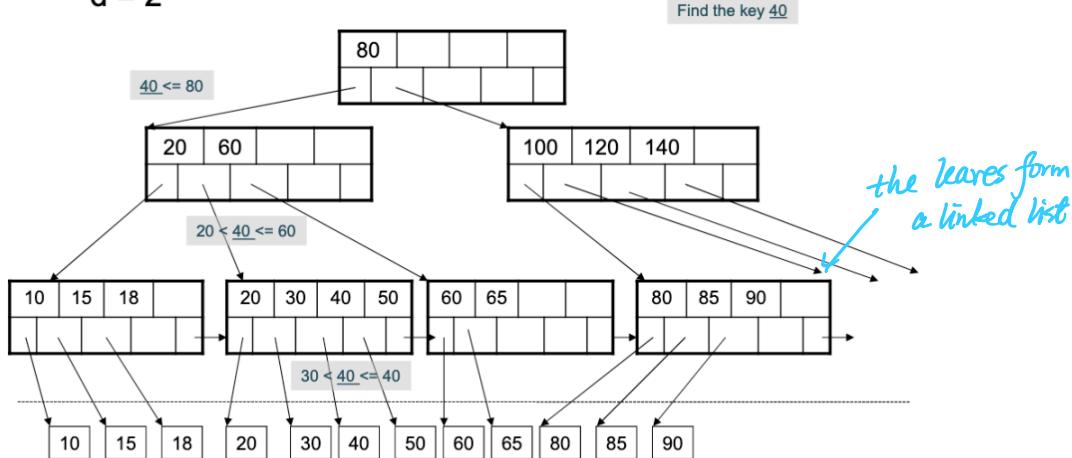
- Index organization

- Hash table

- $B^+$  trees max number of branch is designed on node
      - Every node can have more than one branch
      - Search tree but not binary instead have higher fan-out  
leaves form a linked list

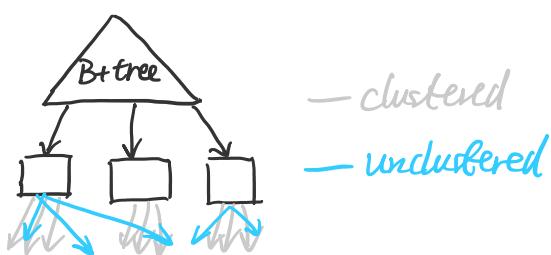
## B+ Tree Index by Example

$d = 2$



- Specialized indices: bitmaps, R-trees, inverted index

- Clustered / Unclustered





- the key of data file must be the same as index.
- We can only have 1 clustered index. → file can only be sorted once and a lot of unclustered indices.
- Clustered: records close in index are close in data
- Unclustered records closed in index may be far in data
- Primary / Secondary

Meaning I: P → primary key. S: otherwise.

Meaning II: P → clustered Secondary: unclustered

## Selection

- "on the fly"

- Index-based selection

logical plan:

$\sigma_{300}$   
Takes

SELECT \*  
FROM Takes y  
WHERE y.courseID=300;

Physical Plan 2:  
index-based

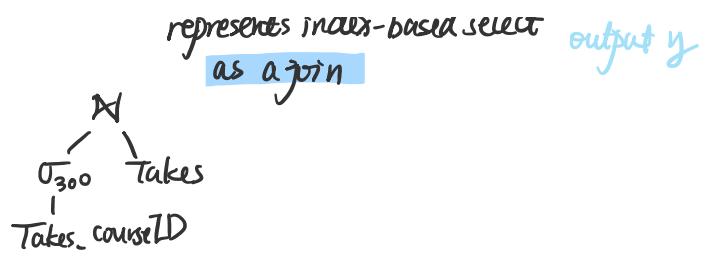
$\sigma_{300}$   
Takes

Take-courseID=  
index on Takes.courseID  
rid-list = Take-courseID.get(300)  
for r in rid-list  
y = student.getRecord(rid)

Physical Plan 1:

on the fly  
 $\sigma_{300}$   
Takes  
for y in Takes  
if courseID=300 then  
output y

SQL Server



- Join Nested-Loop join vs Index-based join

`CREATE TABLE V(M int, N text, P int);`

`CREATE INDEX V1 ON V(N);`

`CREATE INDEX V2 ON V(P, M);` *order matters*

`CREATE UNIQUE INDEX V4 ON V(N)` *the attributes are unique*

**CLUSTERED** *(Not supported by SQLite)*  
*generate a clustered file*

- If there're no updates, just create an index for all attributes.

For range queries, use BTREE. e.g.  $N > 15 \text{ and } N < 15000$ .

e.g. `select *  
from V  
where N = ?  
100000 queries`

`select *  
from V  
where P = ?  
100 queries`

`insert into V  
values (?, ?, ?)`

200 queries

*if there are updates, only create index on N.*

*if not, create index on both N and P.*

*if there're 10000 queries,*

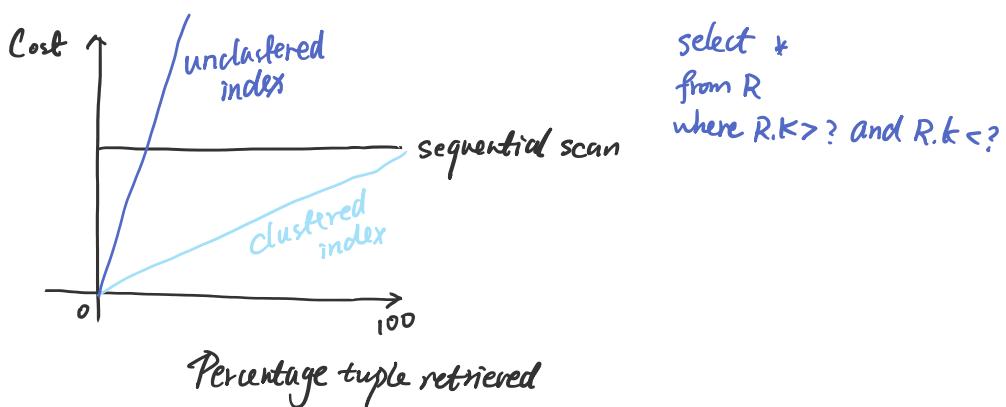
*try using V(N) primary index and V(P) secondary index*

- Two typical kinds of queries

o Point queries  $P = ?$

o Range queries  $P > ?$  and  $P < ?$

- Flash or  $B^+$  tree index
- $B^+$  tree index
- Clustered or not
- Clustered
- Range queries benefit mostly from clustering because they may read more than 1-2%.



- Summary for physical plan:

Access path selection    scan the relation vs index

Implementation choice    (for operator)

Scheduling decisions    pipelined vs intermediate materialization

things go materialization always go blocked.  
 ⇒ they are materialized into an intermediate table.  
 e.g. group by.