

22/01/18

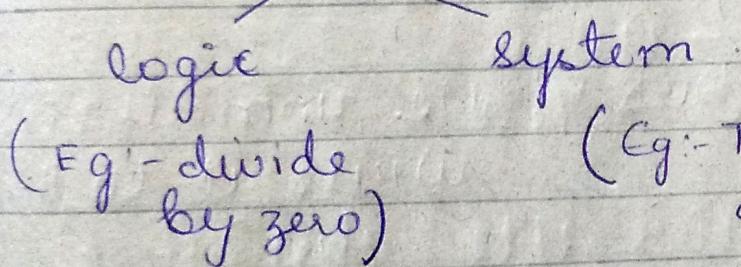
## Q. Recovery

Two properties :-

- ① atomicity.
- ② Durability

Types of failures :

① Transaction failure



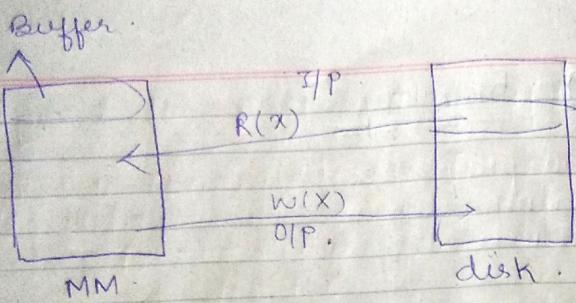
② Disk failure

③ Raid (Redundant array of indepe-  
ndent disk).

→ To recover from natural disaster

To recover data from natural  
disaster, we should store it on different  
geographic location.

④ Crash failure



Work area or work space is known as buffer.

$W(x)$  → Disk (We immediately write it to the disk).  
It is known as immediate mode by  
 $Output(x)$   
commit / abort

If we write to the disk after commit or abort is known as deferred modification (delay).

If failure occurs somewhere b/w  $W(x)$  and  $Output(x)$  → it is known as failure and here we need failure mechanism.

### Log Records

$\langle T_i \text{ start} \rangle$

$\langle T_i \text{ update} \rangle$

$\langle T_i \ X_i \ V_1 \ V_2 \rangle$  (changes  $X_i$  value from  $V_1$  to  $V_2$ )

$\langle T_i \text{ commit} \rangle$

$\langle T_i \text{ abort} \rangle$

Undo → we do undo some  $T_i$  then it set the previous value / old value

Redo → it set the new value.

Eg:-

$\langle T_0 \text{ start} \rangle$   
 $\langle T_0 \ A \ 1000 \ 950 \rangle$   
 $\langle T_0 \ B \ 2000 \ 2050 \rangle$   
 $\langle T_0 \text{ commit} \rangle$   
 $\langle T_1 \text{ start} \rangle$   
 $\langle T_1 \ C \ 700 \ 600 \rangle$   
 $\langle T_1 \text{ commit} \rangle$

crash occurs somewhere here

I Before commit / redo, the value will be undo.  
 $A = 1000$        $B = 2000$

(ii)

Suppose failure occurs after update of C, so T<sub>0</sub> will be committed.

As T<sub>0</sub> is committed, so T<sub>0</sub> will be redo.

so T<sub>0</sub> — Redo.

$$A = 950$$

$$B = 2050$$

As T<sub>1</sub> is not committed yet, so we will undo T<sub>1</sub>.

T<sub>1</sub> — Undo.

$$C = 700$$

If crash occurs after commit of T<sub>1</sub>, so both T<sub>0</sub> and T<sub>1</sub> will be redo.

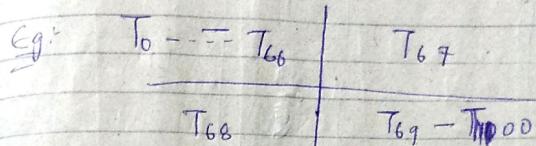
$$A = 950; B = 2050, C = 600$$

### Checkpoints:

Checkpoints will forcefully write the values to disk whether transaction is committed or not and we will add a checkpoint record <checkpoint>

Now we will check the active transaction if crash occurs (after the checkpoint).

RTA) ~~-----~~ (RTB)



T<sub>0</sub> - T<sub>66</sub> committed before checkpoint  
T<sub>68</sub>      "      "

so T<sub>67</sub> and T<sub>69</sub> - T<sub>1000</sub> (are in active state)

### Series of steps:

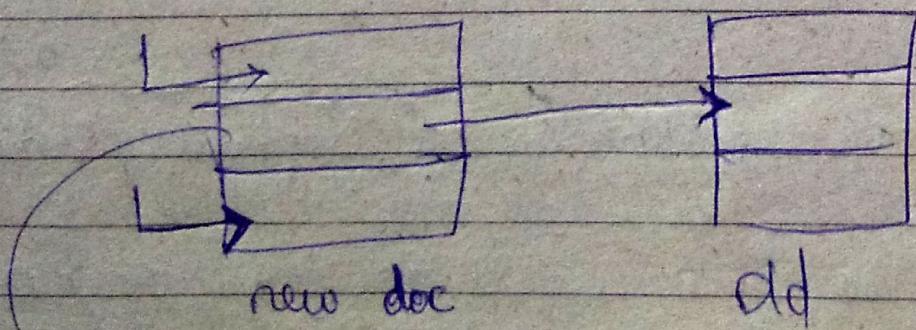
- ① Suspend all the active transaction
- ② Add checkpoint record and maintain a log file
- ③ Forcefully write to disk
- ④ Resume all the transactions

## Shadow Paging

We maintain the copy of the whole database we will make changes to the new database, and we will see if all is written to disk. Then we will make the new database the current database and delete the old database.

~~Copy~~ we can't use this for very large database. So we use paging.

We will divide blocks into pages and maintain a page table. We will access blocks with page number and the page ~~width~~ in which we make modification will be only copied and the one in which there is no modification will point to the old one.



The one with no modification points to the Old one.

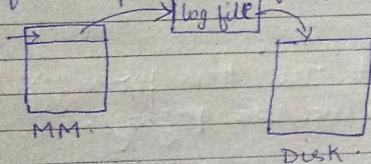
This is known as shadow paging.

But practically, this is also not good for large database.

23/10/18

Algorithms to recover data from log file

① Deferred update (delay update) :-

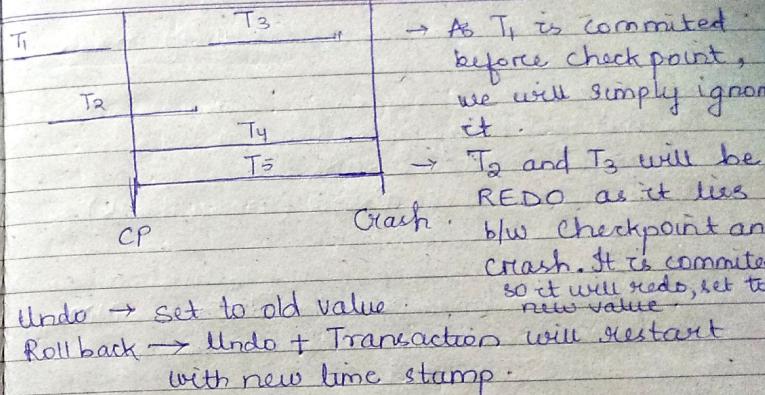
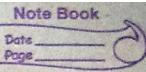
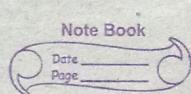


Just Before the commit statement is known as Commit point.

As we reach commit point, then commit statement goes to log file and then all the updates are written to disk.

This algorithm is also known as No Undo/Redo.

Active state of transaction means no abort or no commit.



→  $T_4$  and  $T_5$  will be roll-backed as it is not committed.

Whenever crash occurs, we check from last statement before crash to the first check point occurred.

Immediate Update

Immediate

(UNDO / NO REDO)

(As soon as T is updated we will write it to the disk).

(As all the transaction are being updated simultaneously)

Commit

Here we will first write the log file to the disk and then we will write commit statement to the log file (UNDO / REDO).

→ In immediate update cost is very much

→ In the previous example,

$T_1$  will be simply ignored.

$T_2$  and  $T_3$  will be REDO.

$T_4$  and  $T_5$  will be UNDO (no rollback).

→ We can't store/maintain the ~~data~~ in MM  
as it is volatile, when crash will occur  
log file will be deleted.

so we store log file in disk.

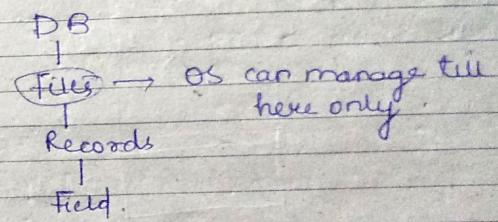
(Ques: But if crash occurs in disk then where  
the log file goes?)

→ If crash occurs while recovering mechanism  
(i.e., after one crash) then what happens?  
(We understand that log file also maintains  
UNDO AND REDO operation).

→ Why do we need DBMS control manager to  
maintain DB files? Why can't OS deal  
with DB files also?

In OS if want to access a file, we can  
access a complete file (not a part of it).

But in DBMS, we can access one record;  
multiple records (rows), one field, etc.  
i.e., we can access a part of file. So OS  
can't deal with it.



26/10/18

### Indexing and File Organization

Database → collection of files

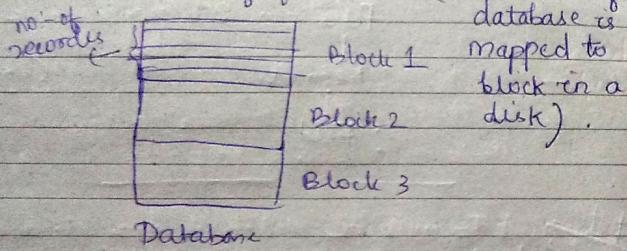
files → Records + fields



Average no:- of records in one block → Blocking factor (Bfr).

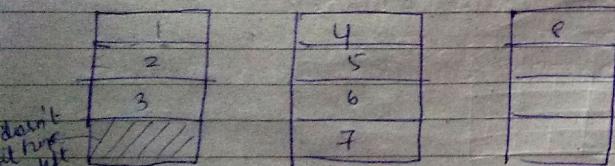
$$BFR = \frac{\text{Block size}}{\text{Record size}}$$

We divide database into blocks and each block consist of records. (We will not talk now in terms of files). (The block of database is mapped to block in a disk).



We can store records in block in 2 ways

### ① Unspanned.



Eg:- of Unspanned.

Unspanned - If our record doesn't fit in one block, we will waste memory and move it to another block

But advantage is:- For one record, we have to access one block only.  
This is used for fixed size records.

### ② Spanned.

If our record doesn't fit, then we will store that much part of record which fit in that block and will move the rest to another block.

Advantage → No memory wastage.

Disadvantage → Have to access two block.

This is used for ~~variable~~ records having no fix size.

Eg:-

URL	Content

The content corresponding to each url is variable, so we will use spanned.

→ When we have ordered database (in ascending or descending order (e.g. eid)), so insertion cost is high, deletion is very low.

But when we delete, some memory space remains vacant, then database reorganize it.

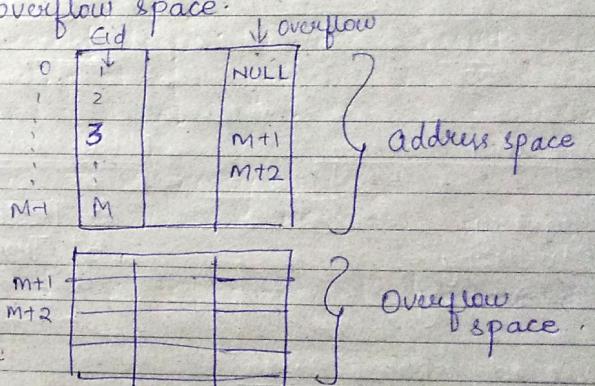
For searching of blocks, we will use binary search  $\rightarrow$  Complexity  $(\log b)$  where  $b$  is no. of blocks.

### Unordered:

- $\rightarrow$  Not ordered in any order.
- $\rightarrow$  Now suppose our database is ordered by eid and we are using ename, then this is also un-ordered.
- $\rightarrow$  Data structure used is heap.
- $\rightarrow$  linear search is used.
- $\rightarrow$  Insertion is easy as it is unordered.
- $\rightarrow$  Here also we do reorganization periodically when we delete records.
- Hashing:
- $\rightarrow$  we have to design hash function in such a manner, if we are applying it on any hash field value then we get the exact index of that field value.
- $\rightarrow$  The O/P of the hash function is the index value.
- $\rightarrow$  Hashing  $\rightarrow$  constant time complexity as it gives exact memory location.

Ordered  $\rightarrow$  can be spanned or unspanned (?)  
and Unordered

$\rightarrow$  Disadvantage  $\rightarrow$  when at same space two records exist  $\rightarrow$  collision; then we do chaining.  
For this we maintain a overflow pointer and overflow space.

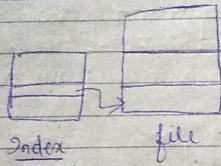


Suppose at Eid 3, we have two records, then we will write the index where the second record is stored in the overflow space, like here it is M+1.

- $\rightarrow$  External Hashing  $\rightarrow$  Apply hashing on blocks (Searching blocks).
- $\rightarrow$  Internal hashing  $\rightarrow$  Apply hashing on records (Searching records in a block).

30/10/18

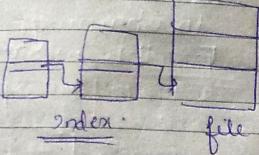
### Single level Index.



### Note Book

Date \_\_\_\_\_  
Page \_\_\_\_\_

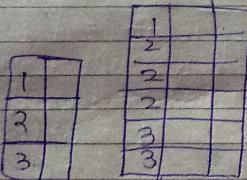
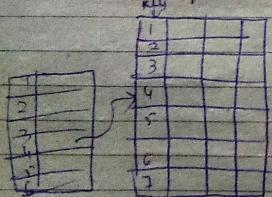
### Multi-level Index



### Single Level Indexing

Primary  
Clustering  
Secondary

### Dense and Sparse Index



Index file

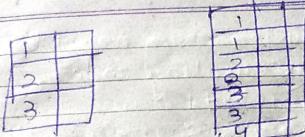
Dense (every key value) <sup>record</sup>  
{ If for every value, there is an index then it is dense }

### Sparse.

(As for every record there is no indexing)

### Note Book

Date \_\_\_\_\_  
Page \_\_\_\_\_

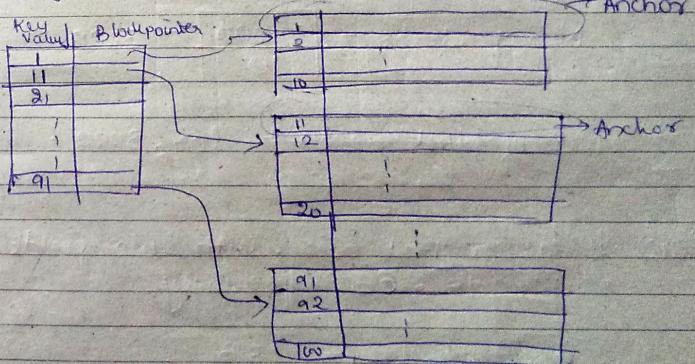


This is both dense and sparse. As for every key value there is an index so it is dense and for every record, there is no index so it is sparse.

### Primary Indexing ↓ Sparse

#### Adv

→ Ordered (and ordered on primary key)  
→ length of the record is fixed.



→ Anchor → The first record of every block.

→ Index are always ordered.

→ This is sparse.

If we are directly accessing block  $\rightarrow$  no. of block access =  $\log_2 b$ ; But by creating index, no. of block access =  $\log_2 b + 1$

- Q There is an ordered file with 30,000 records, stored on a disk with block size 1024 bytes. File are of fixed size and are unspanned with record size = 100 bytes. Find the number of block access.

$$bf = \left\lceil \frac{1024}{100} \right\rceil_{\text{floor}} = 10.$$

$$\text{No. of blocks} = \left\lceil \frac{30000}{100} \right\rceil_{\text{ceil}} = 3000.$$

$$\text{No. of block access} = \log_2 3000.$$

- (ii) If given that index, with two fields each of 9B, 6B respectively

$$\text{Record size} = 15B$$

$$bf = \left\lceil \frac{1024}{15} \right\rceil_{\text{ceil}} = 68 \text{ records/Block}$$

So no. of indexes in index table

$$= \left[ \begin{array}{c} 30,000 \\ 68 \\ 2720 \\ 28042 \end{array} \right] \\ = \underline{45}$$

$$\text{No. of block access} = \underline{\log_2 45 + 1}$$

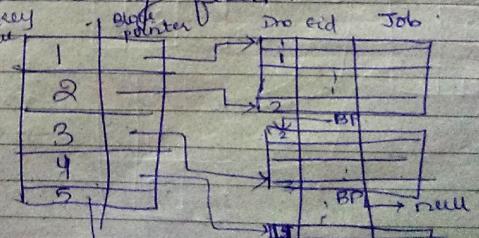
Disadvantage of primary

→ It can only be applied on ordered ~~primary~~ key field and not applicable to unordered field

Clustering Index

ordered + nonkey

We will apply the search on ordered non-key field which is known as clustering field.



→ In index, we will do ~~every~~ entry of every unique record.  
 → It is both dense and sparse.  
 → Here there is a BP (Block pointer) in the file pointing to the next block if the non-key value is same otherwise NULL.  
 → The index will access the occurrence of non-key value in the block. The index will point to the first occurrence of 1 in file.  
~~→ No. of file access  $\geq \log_2 + 1$~~   
Secondary Index:  
 → It is unordered; candidate key (unique values but not ordered)  
 → It can be on Key or non-key  
 → Ordered or primary key.  
 → There will be primary index on A here in this case.

key	Block Pointer
3	B
5	B
9	B
10	C
21	C

Secondary Index:  
~~No. of secondary index nodes = no. of B values~~

	B	C
9	3	1
5	1	10
21	3	

This is the case of when we don't have duplicate B values (non-key).

→ If we have duplicates, we will create a block for every duplicate value and to refer Note Book block we will make another block with ordered & unique value.

Now, for duplicates there will be use of primary indexing (doubt).

No. of records = 30,000  
 Record size = 100 Byte  
 Block size = 1024 Byte

$$Nod\cdot Bf = \left\lceil \frac{1024}{100} \right\rceil = 10$$

$$Nod\cdot of blocks = \frac{30000}{100} = 3000$$

As it is unorderd, so we will apply linear search.  
 So avg no. of block access =  $\frac{3000}{2} = 1500$ .

Now,  
 If we created a secondary blocks with 9B, 6B  
 so one record size = 15B.

$$Bf = \left\lceil \frac{1024}{15} \right\rceil = 68$$

No. of records in index file = 30,000 (or entries)

No. of blocks in secondary index file =  $\frac{30000}{68} = 442$ .

Y18 + 22 + 10

Note Book

Date \_\_\_\_\_  
Page \_\_\_\_\_

\* Index is always ordered, so we can apply binary search  
so no. of searches =  $\log_2 448 + 1$

31/10/18

Record	int eno	2
	char name [22]	22
	float sal	4
	char dep [10]	10
		38 = Record size

No. of Record = 5000.

Block size = 512 bytes.

① No. of block access without indexing.

②

$$Bfr = \left\lceil \frac{5000}{512} \right\rceil = 13 \text{ records.}$$

$$\text{No. of blocks} = \left\lceil \frac{5000}{512} \right\rceil = 385 \text{ blocks.}$$

For ordered, binary search,

$$\text{No. of block access} = \log_2 (385).$$

② Binary indexing (15 bytes each)

$$\frac{385}{15} = 34 \text{ record index.}$$

01/11/18

indext  
No. of blocks =  $\frac{5000}{34}$   
= 385  
34

$$= 12.$$

$$\text{No. of searches} = \log_2 12 + 1$$

01/11/18

### Multi-level Indexing

→ We do levelling till our all index comes in one block.

→ In that indexing which ~~with~~ all the index comes is known as top level index.

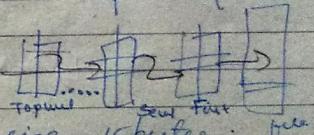
→ We do this to reduce the number of block access.

Q. No. of rec 30,000.

block size - 1024 bytes

record size - 15 bytes

secondary index record size - 15 bytes



$$Bfr = \left\lceil \frac{1024}{15} \right\rceil = 68. \quad (\text{Because here it is already given that secondary index is there so we need not divide 1024 by 15})$$

$$\text{no. of entries} = 30,000.$$

$$\text{no. of blocks}, \left\lceil \frac{30,000}{68} \right\rceil$$

4.1.2.2 + 10

Note Book

Date \_\_\_\_\_

Page \_\_\_\_\_

As Index is always ordered, so we can apply binary search  
 No. of searches =  $\log_2 448 + 1$

31/10/18

Record	int eno	2
	char name [22]	22
	float sal	4
	char dep [10]	10
	38	= Record size

No. of records = 5000.

Block size: 512 bytes.

① No. of block access without indexing.

$$\text{Bfr} = \left\lceil \frac{5000}{512} \right\rceil = 13 \text{ records.}$$

$$\text{No. of blocks} = \left\lceil \frac{5000}{512} \right\rceil = 385 \text{ blocks.}$$

For ordered, binary search.

$$\text{No. of block access} = \log (385).$$

② Primary Indexing (15 bytes each)

$$\left\lceil \frac{385}{15} \right\rceil = 26 \text{ record/index}$$

$$\text{No. of blocks} = \left\lceil \frac{5000}{385} \right\rceil$$

Note Book

Date \_\_\_\_\_

Page \_\_\_\_\_

$$\begin{array}{r} 385 \\ \hline 34 \end{array}$$

$$= 12.$$

$$\text{No. of searches} = \log_2 12 + 1$$

31/10/18

### Multi-level Indexing

→ We do levelling till our all index comes in one block.

→ In that indexing which ~~with~~ all the index comes is known as top level index.

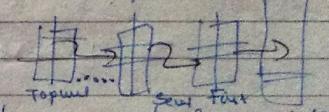
→ We do this to reduce the number of block access.

Q. No. of records 30,000.

Block size - 1024 bytes

Record size - 15 bytes

Secondary index record size - 15 bytes



Top level secondary index

$$\text{Bfr} = \left\lceil \frac{30000}{1024} \right\rceil = 68. \quad (\text{Because here it is already given that secondary index is there so we need not})$$

$$\text{No. of entries} = 30,000.$$

$$\text{No. of blocks, } \left\lceil \frac{30000}{15} \right\rceil$$

$$= 442$$

As, 442 is also more blocks.  
so, we will do <sup>2nd level</sup> indexing

$$\text{So no:- no:- of blocks} = \left\lceil \frac{442}{68} \right\rceil = 7$$

Now 7 is also more than 1.

so we will go for 3rd level indexing.

$$\text{So no:- of blocks} = \left\lceil \frac{7}{68} \right\rceil = 1$$

So now here, 3rd level index is our top level index as here there is only 1 block.

In the original file ; no:- of records per block

$$\text{no:- of blocks} = \frac{30000}{10} = 3000 \quad \begin{matrix} 10 \\ \text{blocks} \end{matrix}$$

So from 3000 blocks  $\rightarrow$  442 blocks  $\rightarrow$  7 blocks  
original file. 1st level and 2nd level

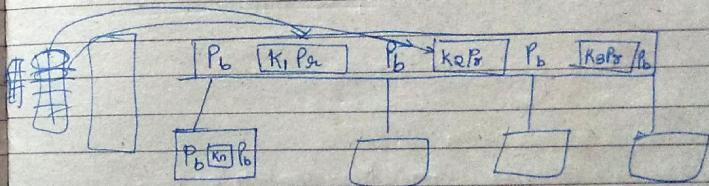
1 block  
top level

$$\text{No:- of block access} = 4. (1+1+1+1)$$

$$\text{No:- of block access} = t + 1 \quad \begin{matrix} \xrightarrow{\text{file access}} \\ (\text{top level}) \end{matrix}$$

$\rightarrow$  Here we will have problem for insertion & deletion. So to overcome this we move towards dynamic multi-level indexing which is implemented through B and B<sup>+</sup> tree.

B-tree  $\rightarrow$  structure of internal node & root node is same.

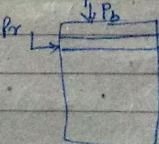


Pb  $\rightarrow$  Pointer pointing to a particular block

Px  $\rightarrow$  Pointer pointing to a particular record

No:- of blocks < No:- of records.

Size of block pointer < Size of record pointer



Order of the B tree = maximum no. of children  
 Minimum no. of ~~children~~ children = 2

of root  
 $\min = p$

For root, key  
 Value can be -  
 $\min = 1$   
 $\max = p-1$

node { No. of children, min =  $p/2$ , max =  $p$ .  
 No. of key values, min =  $\lceil p/2 - 1 \rceil$ , max =  $p-1$

Q Consider a B tree with key size 10 bytes,  
 Block size = 512 bytes, data pointer ( $P_d$ ) size  
 = 8 bytes and Block pointer size = 5 byte.  
 Order of B tree.

No. of block pointers =

$$nP_b + (n-1)k + (n-1)p_r \leq \text{Block size}$$

$$\Rightarrow n \times 5 + (n-1)(10+8) \leq 512$$

$$\Rightarrow 5n + 18n - 18 \leq 512$$

$$\Rightarrow 23n \leq 530 \Rightarrow n \leq 530/23$$

$x \rightarrow y$   $y \leq x$   
 $ycx$

Note Book  
 Date \_\_\_\_\_  
 Page \_\_\_\_\_

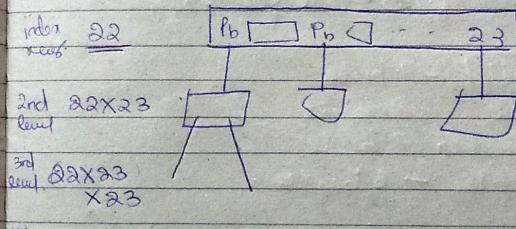
$$\Rightarrow n \leq 23$$

$$\underline{n = 23}$$

02/11/18

Q Order of B trees = 23

Then how many maximum index records can be stored in 4 levels of B tree.

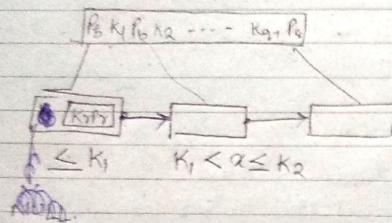


Total index records in 4 levels =

$$22 [ 1 + 23 + 23 \times 23 + 23^3 ]$$

$$= 22 [ 1 + 23^1 + 23^2 + 23^3 ]$$

## B+ trees



Q Find the order of B+ tree for the given data.

Index key field - 9 byte  
Block size 512 byte  
Record pointer 7 byte  
Block pointer - 6 bytes

Internal node

$$n \times 6 + (n-1) \times 9 \leq 512$$

$$\Rightarrow 15n - 9 \leq 512$$

$$\Rightarrow 15n \leq 521$$

$$\Rightarrow n \leq \frac{521}{15}$$

$$\Rightarrow n \leq 34$$

$m = 34$  Order for internal node

Note Book

Date \_\_\_\_\_  
Page \_\_\_\_\_

study views from lab

Note Book

Date \_\_\_\_\_  
Page \_\_\_\_\_

Leaf node

$$(n-1)(k+p_r) + nxP_b \leq BS$$

$$(n-1)(7+9) + n \times P_b$$

$$\Rightarrow (n-1) \times 16 + n \times 6 \leq 512$$

$$\Rightarrow 22n - 16 \leq 512$$

$$\Rightarrow 22n \leq 528$$

$$n \leq \frac{528}{22}$$

$$\leq 24$$

$$n = 24$$

Order for leaf node  
→ pointing to next node

$$\Rightarrow m(k+p_r) + P_b \leq BS$$

$$\Rightarrow m(7+9) + 6 \leq 512$$

$$\Rightarrow 16m + 6 \leq 512$$

$$\Rightarrow 16m \leq 506$$

$$\Rightarrow m \leq \frac{506}{16}$$

$$\leq 31.6$$

$$m = 31$$

order for leaf node