

AGENDA

- 1.) what is index / why / advantages
- 2.) working of index (B/B+)
- 3.) cons of index
- 4.) index on Multiple cols
- 5.) index on strings

INDEX:

JOINS:

```
→ for i.....M:  
    for j.....N:  
        if cond match:  
            dosomething()
```

✓ slow $\rightarrow O(N \cdot M)$

Assume \rightarrow 1M Rows ...

Query is slow.

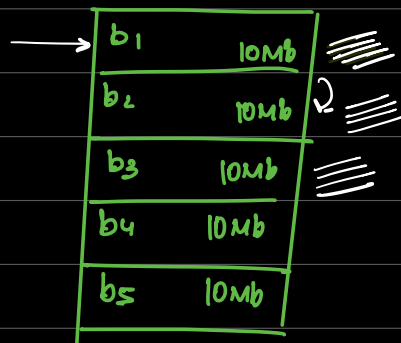
\rightarrow d/b stores data ON DISK.

RAM v/s DISK

✓✓



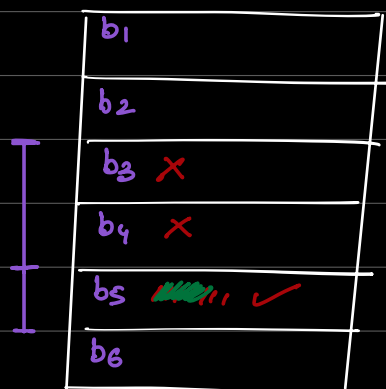
* Internal working of Disk:



(b₁)
 ↓
 from Disk → Memory
 Memory → executed.

PROBLEMS:

✓ select * from students where i = 100;



d/b → b₃/b₄/b₅

① block contains data.

(Disk)



Presence of idx helps you to speed up work!!

"Indexes sort the data" ——— xxxxx

Purpose → to reduce No. of Disk Access.

★) WORKING OF INDEX:

Assume → 100M Rows (table)

↓

Query → select * from stud where id = 100;

Yes
Pr

id	Name	email
1	A	...
2	B	...
3	C	...
4	D	...

(100M) Rows

→ use idx index

DS used ↓

<K,V>

<id, address_block>

<K> <V>

<u>id</u>	value
1	#abc12
2	#abc123
3	#q1172
⋮	

Address
of Memory block.

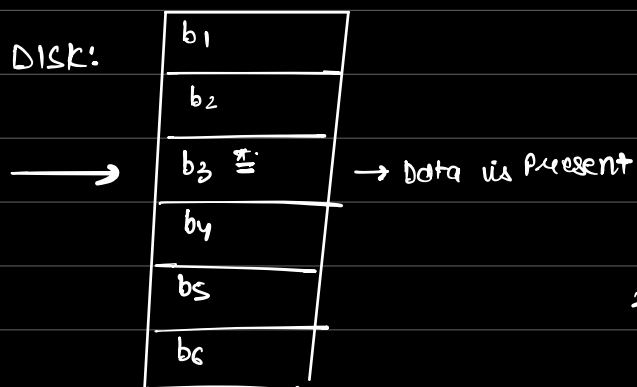
100 ↓ #aa100bc
.....

(stored in memory.)

select * from stud where id = 100;

s1.) before Query execution → Go to Map, Get address of block.

s2.) d/b Gets data from block.



1.) No Index : 3 blocks ↓

2.) Index table : 1 block ✓

'id' → <k> index table.

↓↓↓

id	Name	email_id	...
----	------	----------	-----

Name →

" select * from stud where name = 'yash' ;

10 yash

	key	↓ value
xxx	yash	#1, #2, #3

value : List<string> = address blocks.

<Yash, [odd1, odd2, odd3 ...]>

→ still Not fetching all blocks

<key, C> → work for every Query??

3.2 Range based Queries

Index table

id	Name	psp ^(*) → index
1	A	10.0
2	B	10.1
3	C	10.2
4	D	10.3
		⋮

10.0	b1
15.2	b2
21.7	b3
70.1	b3

select * from stud where
psp > 10 and psp < 60;

}

for i=10 ... 60:

// go to index table

⋮

10-60

10.6

10.6

10.0

10.1

10.2

10.3

⋮

10.7

11.0

this DS: $\langle \text{key}, [\text{values}] \rangle$

x won't work.

**) REASONS:

- HashMap based indexes won't work.
- fail on Range based queries.

solutions:

Requirements for Index:

-) data to be ordered. → (Index table)
i.e. (key)
- value → block-address.

$\langle k, v \rangle$
↓
(sorted)

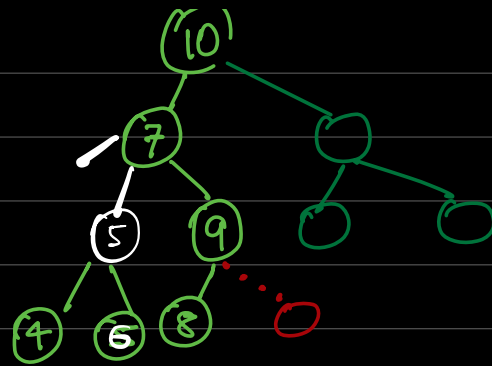
Treeset / ordered-Map → solve

- ✓ keys sorted
- ✓ stores $\langle k, v \rangle$

TC: $O(\log N)$

$O(1)$ v/s $O(\log N)$

Internally they are BST. (balanced bst)



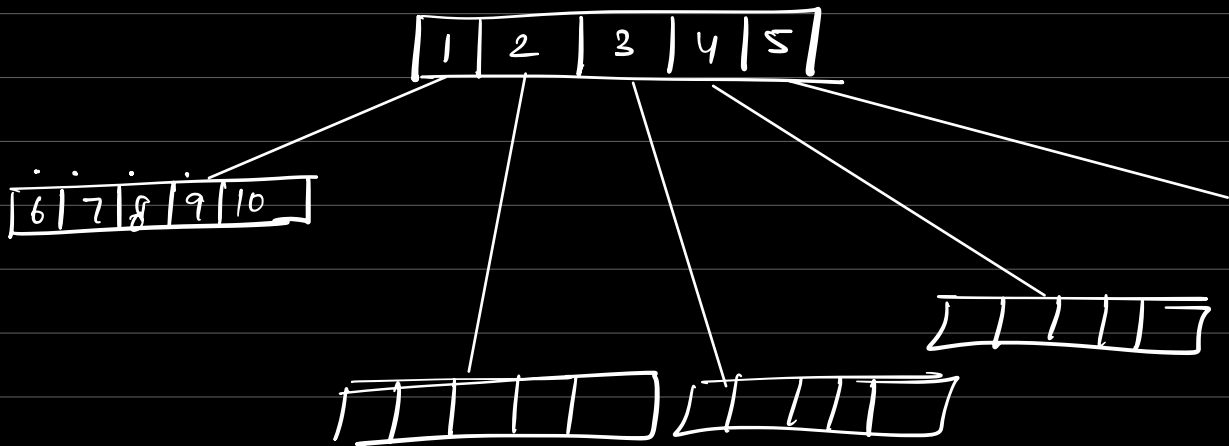
7-10

$\log(N)$
(height)

Indexes \rightarrow Internally work ON B. BST.

TreeMap \rightarrow Internally uses B/B⁺ trees

exact B. BST ; Just ①
difference: $\leq x$ Nodes.



B/B⁺ \rightarrow B. BST \oplus every Node can have $\leq x$ Nodes.
here $x=5$

-) How index work
-) benefits
-) why not maps x

10:29 PM IST

★) CONS:

when filled ??

Insert
delete
update } → write queries

Read → Read query.

1.) write queries are slower

2.) Index table is stored on RAM & copy is stored on disk.
↓ Extra space.

Storage increases.

Main point → Not always use indexing.

*) INDEXING ON MULTIPLE COLS°

id ✓

(id, name) ✓

(id, name, psp) ✓

99% make mistakes.

" select * from stud where psp = 20; "

idx ON: (name, psp)

→ (X) NOT be faster (NOT used)

created a index ON (NAME)

select * from stud where name = "Yash";

•) sorting names

•) storing index table

	AMAN	#1
	Payal	#2
→	<u>Yash</u>	<u>#3</u>

v/s

$(name, psp) \neq name \text{ and } psp$ xxxxxx

Group by b-id, psp

↓
??

Similarly in Index $psp \rightarrow$ acts as tiebreaker

Name psp

Aman 70

Aman 90

Aman 80

Naman 20

Yash 70

Index ON (Name, psp)

→

⇒

Aman 70

Aman 80

Aman 90

Naman 20

Yash 70

Index (Name, psp) = Index (psp, Name)

↑↑

↑↑↑

↓

⇒

20

→ 70

70

80

90

Naman

Aman

Yash

Aman

Aman

Query ON	Index-ON	WORK / NOT WORK
① Name	PSP	x
② Name	Name	✓
③ Name	(PSP, Name)	x
④ Name	(Name, PSP)	✓
⑤ $\text{name} = x$ AND <u>$\text{PSP} = 10$</u>	(<u>PSP</u> , Name) ↑↑	~
⑥ $\text{name} = x$ OR <u>$\text{PSP} = 10$</u>	(PSP, Name) ↑↑↑	xx

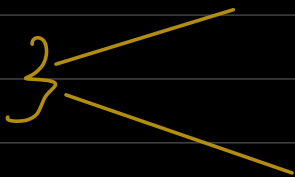
Rule of thumb →

Query on col = x → fast

if Index starts with (x)

- idx (Name)

- idx (PSP)



*] INDEX ON STRING:

users
id email phone

Most common \rightarrow select * from users where
email_id = '.....';

email is indexed \checkmark

<u>email-id</u>	<u>address</u>
abc@scaler.com	#1
yashJain super cool@scaler.com	#1
yash199@scaler.com	#2

users = 1B \Rightarrow 10^9

1.) space

2.) search on string.... x

1B $\rightarrow 10^9$

1B Records

40 char $\times 4 = 160$ bytes

4 bytes $\times 10^9 =$

$$\underline{11 \text{ char} \times 4 \text{ bytes} \times 10^9 = 44 \times 10^9 = 4.4 \text{ GB}}$$

whenever working with indexes on strings \rightarrow

try to use indexes on 1st 5-6 characters....

Google → 2B users

Indexing on Name

Yash Jain

Yash Jainn

Index on 1st 7-8 charac....

Yash Jainwal

→ Name address

Yash Ja (#1, #2, #100)

→ Aman (#2, #110, #41, #22) → 10 blocks / 10

else:

Name	address	
Yash Jain	#1	xxxx
Yash Jainn	#2	too much space
Yash Jainwal	#3	

a-z

10000

Index - ON - chare	users count	# distinct keys	# disk access
1	2×10^9	<u>26</u>	$10^9 / 26$
2	2×10^9	<u>26</u> <u>26</u>	$10^9 / 26 \times 26$
3			
4		---	$10^9 / 26^4$

Query \rightarrow " select * from users
where name like '%yash%';

(full text index.)