

INDEXES

Agenda

- ① What are indexes
- ② How indexes work
→ B and B+ tree
- ③ Cons of indexes
- ④ Indexes on multiple columns (How they help)
- ⑤ Indexes on strings
- ⑥ Indexes PRACTICAL
↳ EXPLAIN.

Till now we majorly focused on writing queries or how do they function via pseudo codes.

We always said pseudo code is not the exact working code, just a medium to build intuition of how things might be happening behind the scenes.

e.g. for joins.

```
for i — 100 :  
    for j — 100:  
        do something(); } } joins
```

if the DB was exactly doing this.

⇒ Performance will be very slow.

① DB have to do a lot of optimizations to make queries fast.

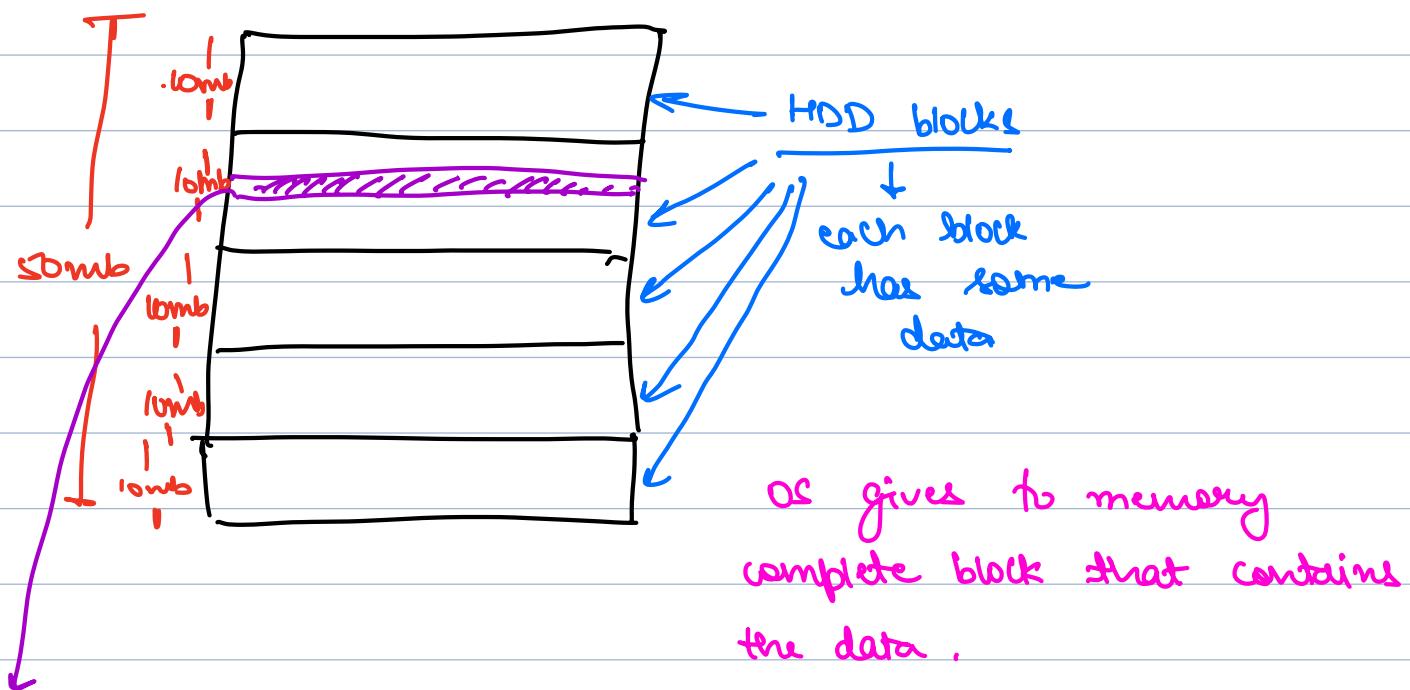
② DB stores data on disk.

vs Program → stores in memory faster

Show table of latency for disk vs memory.

HDD \Rightarrow RAM slower \Rightarrow RAM

Hard disk



Memory (entire block)

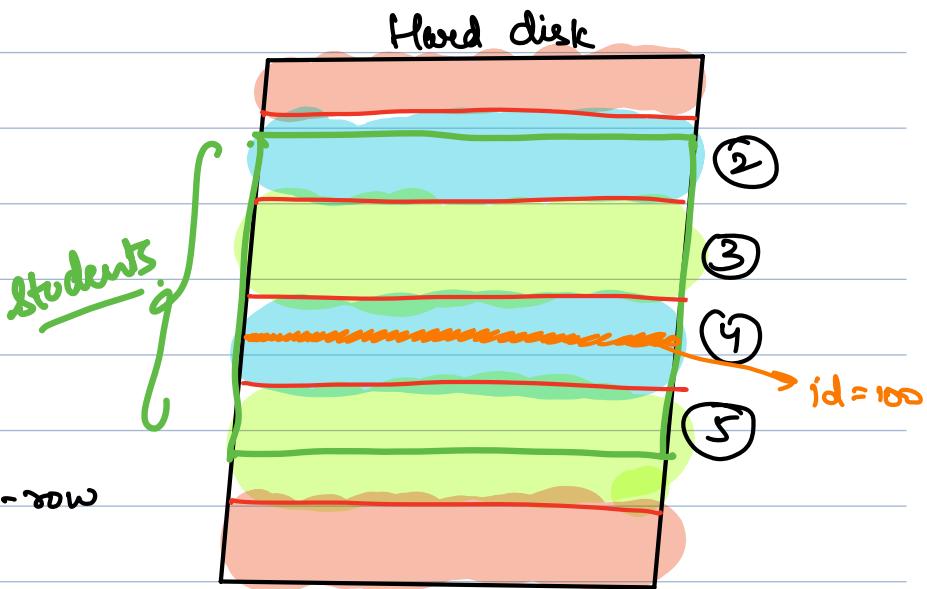


CPU \rightarrow always fetches data from RAM.

If system wants data from disk, it is first loaded into the memory and then read by CPU.

Select *
 from students
 where id = 100

Brute force way for this?
 # Can't traverse row-by-row
 on disk



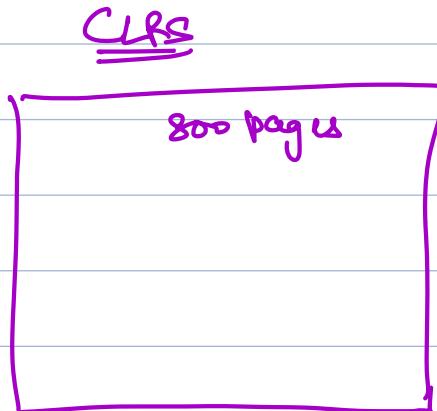
- ① bring block 2 to memory
- ② check all rows if they contain id = 100; next block

NOT OPTIMAL

⇒ unnecessary blocks that are being fetched are a waste of performance.

Where do you see index in day-to-day life?

INDEX		
S.no	Title	P. No



⇒ at the end of a good book, there is an index page, which contains where exactly topics are present.

Index in a book helps find page faster

Index in dB ⇒ finds blocks of disk faster

↓
blocks that contain the rows
that we desire.

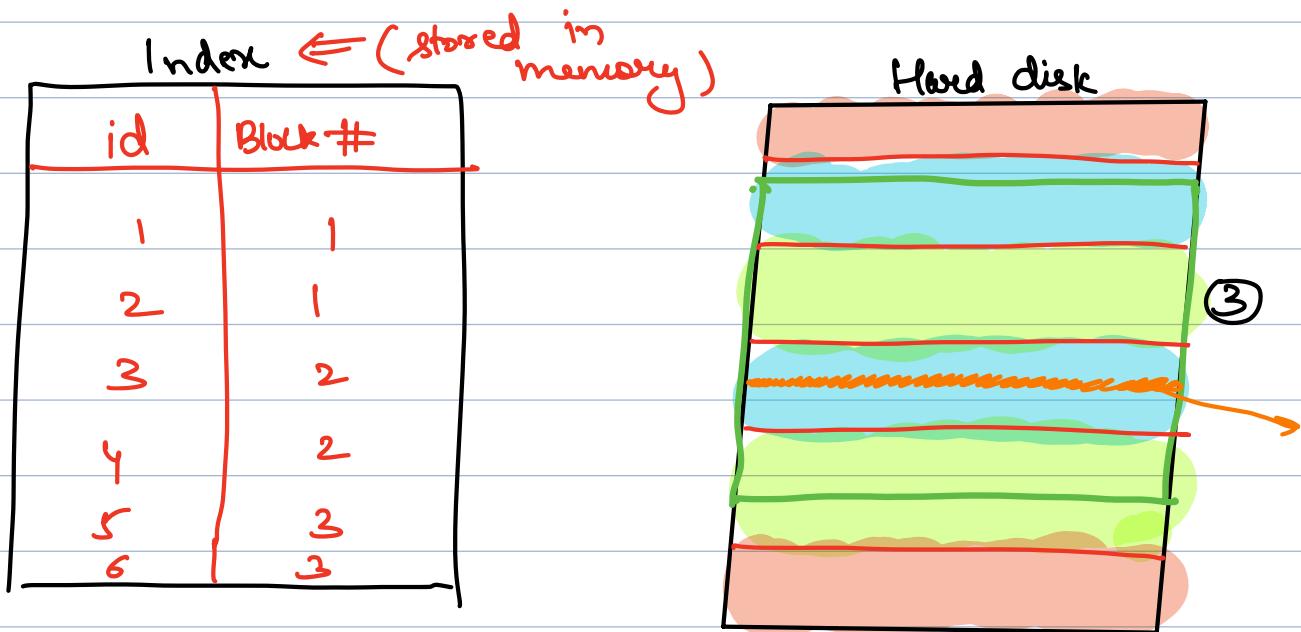
Myth - Indexes don't sort data on the disk

Purpose of indexes : Reduce no. of disk block accesses to fetch data.

HOW INDEXES WORK

Huge table 100M rows.

Select * from students where id = 100 ;



Map /HashMap

database that tells me student with id = 100 is
in this data block directly.

(key | value pair)

without hashmap \rightarrow 4 blocks

with hashmap \rightarrow 1 block

Q Do indexes help with non-unique columns?

Select * from students
where name = 'Ujjwal';

students		
id	name	psp
1	Ujjwal	80
2	Ankit	70
3	Ujjwal	100

Index.

list (block)	
name	block #
Ujjwal	(1, 2)
Ankit	1
Alok	(1, 3)
Deepak	[2, 4, 5]

Select * from students
where name = 'Ujjwal';

without index \rightarrow 6 blocks

with index \rightarrow 2 blocks

Till now we have done indexes

for queries like

where $x = y$;

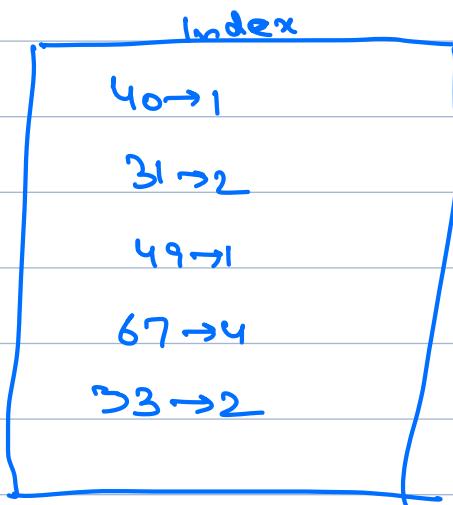
\Rightarrow there are also queries about ranges
eg find all students with name $>= M$
and name $\leq O$;

↓
Starts with M, N, O

e.g. find all students with psp between 60 & 70.

↓

Can we hashmap optimize it.



hashmap < int, string >
get all keys with value
 ≥ 20 and ≤ 40

for i = 20 → 40 :

if i in map:

map.get(i);

since psp is float type
we miss values like

20.1, 31.2 etc.

With a hashmap it is very tough to get all
values within a range.

HashMap → Check a value $\Rightarrow O(1)$

get all values in a range $\rightarrow O(N)$

range queries are like $[L \dots R]$

① Sorted

②

Key : Value

Sorted

Block ID

Map

order-map

Treemap

{

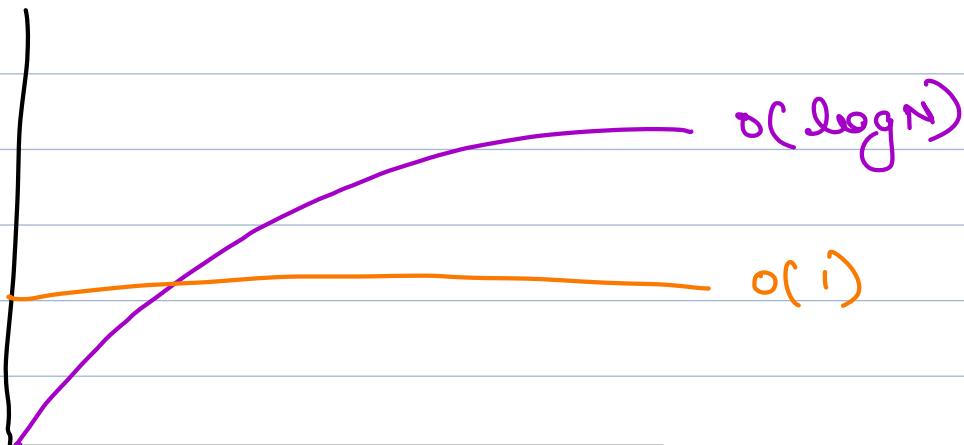
Sorted
map

They work on

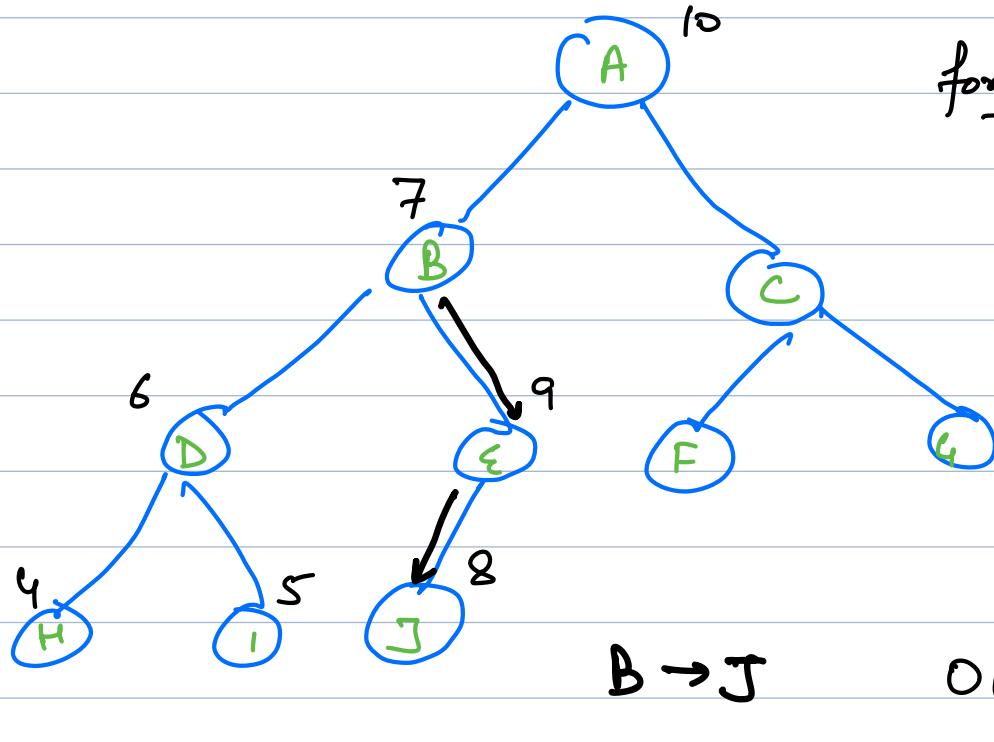
Balanced BST

Check a value $\rightarrow O(\log N)$

difference
isn't
huge



get all values in a range.



for BST.

given a node
get to the
next bigger
node

40 → 60

- ① go to 40
- ② keep going to next node till you go > 60.

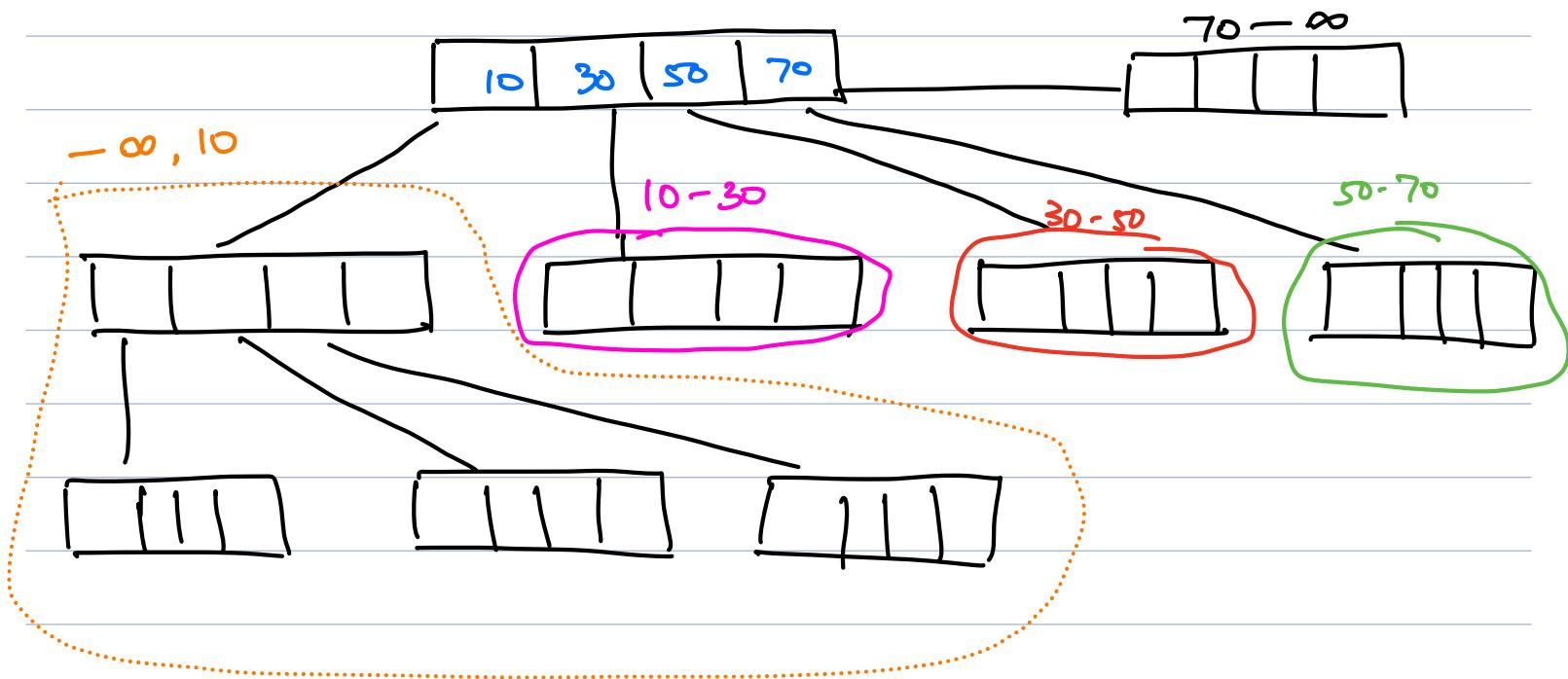
Indexes work very similar to TreeMap.

DS that is used to reduce no. of disk access to increase performance of SQL query.

Tree maps use → AVL or RBT
(eg of BBST)

★ Indexes use data structures called B/B⁺ Trees.

→ instead of every node having ≤ 2 children
every node can have $\leq x$ children.



⇒ Because a B tree has many children
⇒ The height of B tree will be lesser.

⇒ $O(H)$

↑
all of them will be faster for Btree

Checking a value $\rightarrow O(H)$ } faster
getting in a range $\rightarrow O(H)$

What you need to know:-

- ① How indexes work
- ② What is benefit of index
- ③ Why hashmap is not good.

Cons of indexes

When will you put data in a map.

name	blocks

index.

C → ✓

R → X

U → ✓

D → ✓

.

which operations will require
some changes / update in
the index.

Write Operations

(C, U, D) → Change the data
on the disk

- ① writes are slower

we won't keep indexes only in memory.

Typically a copy of index is also stored on disk.

↳ Now for a write you will also need to update this index, making writes even slower

② Storage requirement increases.

* Don't create indexes prematurely.

→ Create an index only when you see the need.



Check performance metrics

⇒ avg time of exec" of queries.

Indexes on multiple columns

(id)

(name)

(name, psb)

(name, psb, rank)



We can create an index on multiple cols as well.

Assume I have an index on (name, psp)

Following query.

Select *

from student

where psp = 80;

Will the above index help? No.

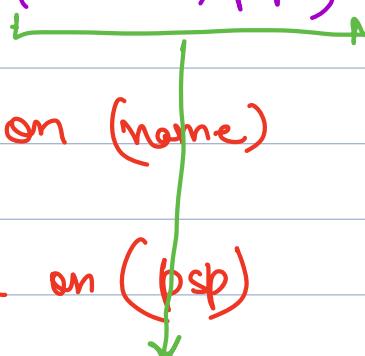
* If you created an index on (name, psp)

#

Creating an index on (name)

and

Creating an index on (psp)



You sort by 2 columns
means

Sorted by name

↓ psp → only tie breaker

Index

(name, psp)

(Abhishek, 70)	
(Abhishek, 90)	
(Abhishek, 100)	
(Bipin, 30)	
(Biswa, 71)	
(Vijwal, 20)	

Sort by name &
Resolve clashes by psp.

(psp, name)

(20, Vijwal)	
(30, Bipin)	
(71, Biswa)	
(70, Abhishek)	
(90, Abhishek)	
(100, Abhishek)	

If you have an index on (name, psp) it won't help a query on only psp.

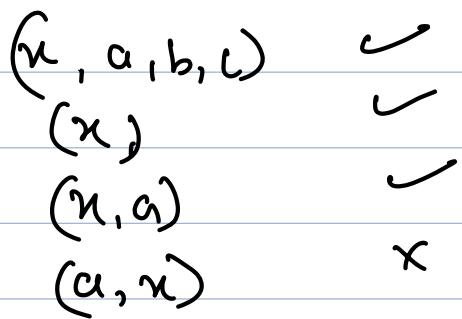
(psp, name)

→ v/s if you have an index on (psp, name) will help.

SCENARIOS.

query on	index on	Help	Not Help
① name	psp		✓
② name	(name) (psp)	✓	
③ name	(psp, name)		✓
④ name	(name, psp)	✓	
⑤ name = x and psp = y	(psp, name)	✓	
⑥ name = x or psp = y	(psp, name)	✓	
⑦ name = x			

★ If a query is on coln x
any index with the x in the prefix will help.



Scenario

$x = a$ and $y = b$

Indexes.

(x)	(y)
<u>(x, y)</u>	

If this is not there
either of them are
equally same.

\Rightarrow SQL finds index with biggest overlap on prefix

$x = a$ and $y = b$ and $z = c$

(x)
(x, z)
(x, y)
(x, z, y) \leftarrow

* AND query multiple coln index will help.
OR query multiple coln index WONT help.

Indexes on strings

users

id	name	email

```
select *  
from users  
where email = ' '
```

index on email coln

email	block
ujjwal@scaler.com	2
smart guy@gmail	10
ujjindal@gmail	40

- ① Size of index will be huge
- ② String matching will be slow

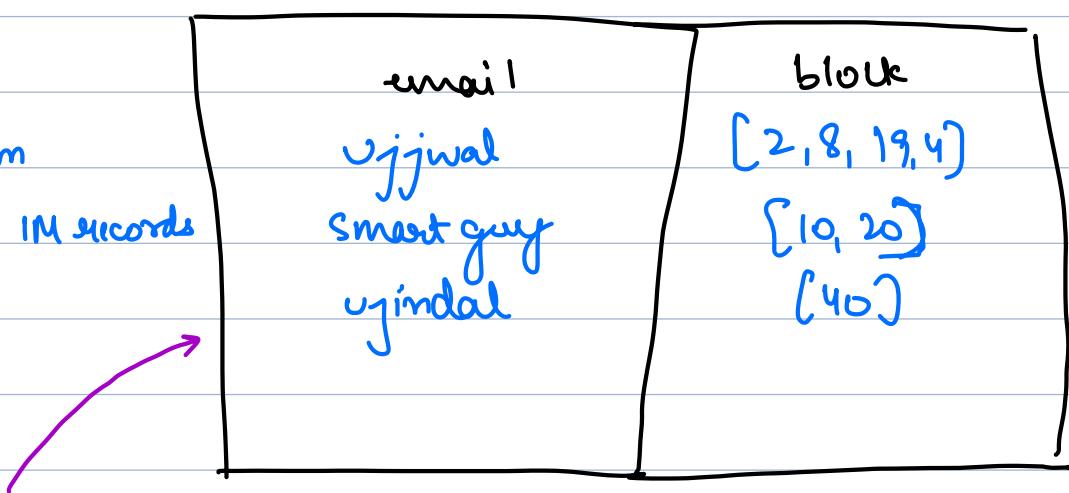
How can we solve this?

Instead of creating an index on complete email,
we create index only on part before '@'.

@yahoo.com

4Bytes x 1M records

40GB



→ Size of the index is reduced.

① Space is saved

Select * from user

where email = ujjwal@scalex.com \Rightarrow will have fetch

4 blocks now as supposed
to 1 block when we
had index on the whole.

without index \rightarrow 250 M

with index on complete \rightarrow 1

with index on first half \rightarrow 4

Trade off \rightarrow Most optimal read or saving 40GB space.

Typically when you have string columns,
index is created on a prefix of
the coln instead of entire column.



Typically that is enough

cg address → index on address [:5]

Countries → countries [:5]

Scalier

⇒ Currently has 2M users / Google 28 users
 2×10^9

index on
how many
characters

users

distinct keys
in index

disk access)

1

2×10^9

26

$2 \times 10^9 / 26$

2

2×10^9

26×26

$2 \times 10^9 / 26 \times 26$



[CODE]

```
for i in range(20):
    print((2 * 10 ** 9) / 26 ** i)
```

if you create an index on the first 7 characters
of address, you will have to access 1 block
on average.

Reduced space required & equally performant.

```
Q select *
from users
where address like '%.ambala%. ;
```



will run index on address column

IT WILL BE OF NO USE.

An index on string doesn't help in pattern matching.



out of scope For this you have
FULL TEXT INDEX

How to create index

good naming practice

[CODE]

Create index idx_film_title_release_year
on film (title, release_year);

explain analyse select * from film

where title = 'ADAPTATION HOLES';



with & without index

Create index idx_film_title on film (title(7));

(slower than complete index)

but still faster than no index)