

Agenda

- Introduction to Indexing → What?
 - How Indexes Work → How?
 - Indexes and Range Queries
 - Data Structures Used for Indexing (Extra material to read)
 - Cons of Indexing
 - Indexes on Multiple Columns
 - Indexing on Strings
 - How to Create an index
-

→ Views not done yet

Quick Introduction

- Graduated in 2019 from Thapar University
- Samsung Research in Noida
- Razorpay (Backend Engineer) (Bangalore)
- Booking.com (Amsterdam) (SDE 2)

Introduction to Indexing

$$10^9 + 10^9 = \underline{\underline{10^{18}}}$$

```
for i .. N:  
  for j .. 100:
```

```
    doSomething();
```

{

Joins

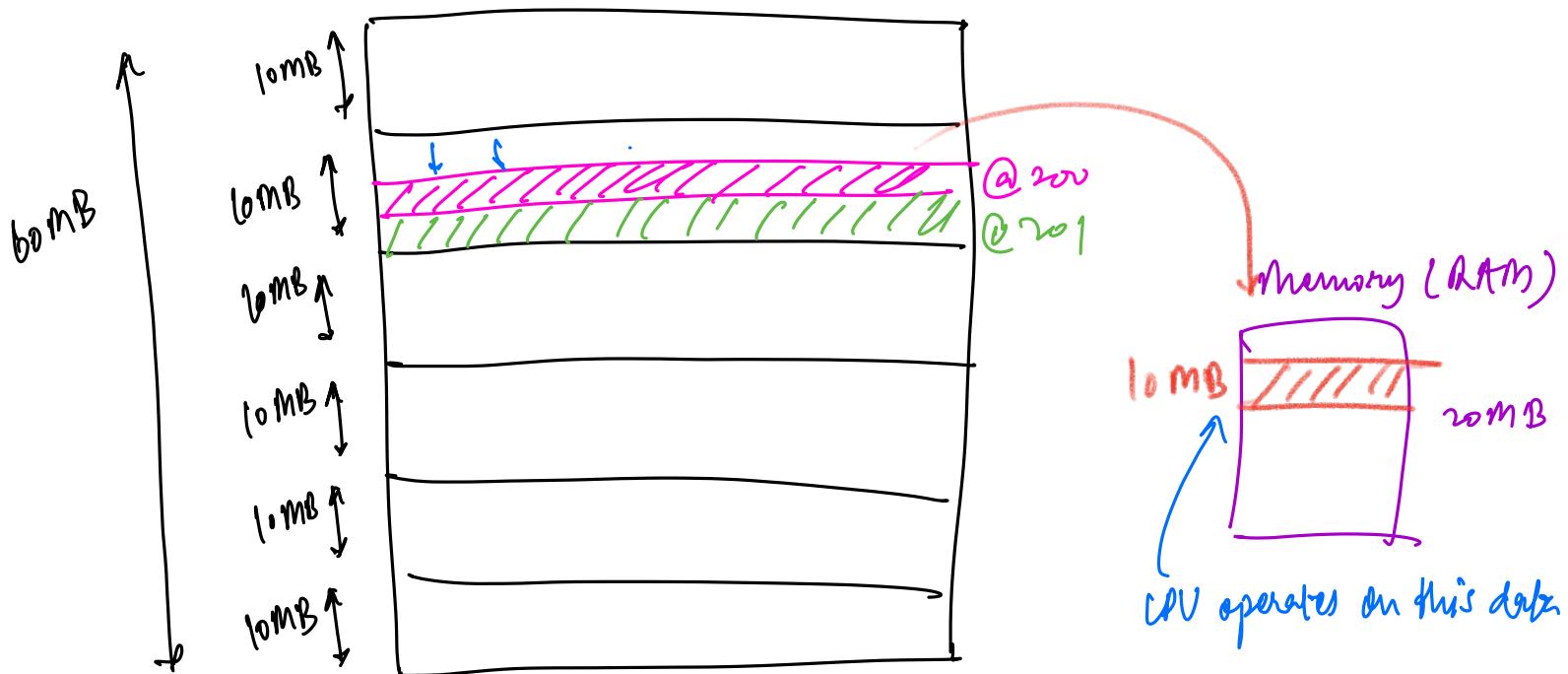
 $\Theta(N^2)$

→ performance will be very very slow

① We need to do some sort of optimization to make our queries fast.

② DB stores data on disk. v/I Program \Rightarrow RAM (very fast)
(slow) Cache (4MB)

HDD is 80x slower than RAM (memory)

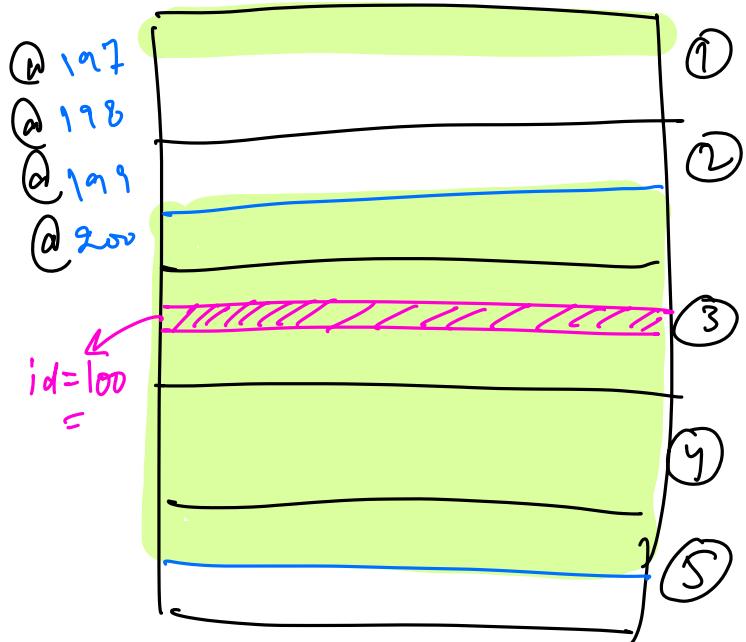


Imagine a table with 100M rows.

Select * from students where $id = 100$;

① Bring blocks ~~2, 3, 4, 5~~ to memory

② Check all the rows for student-id = 100.



→ Unnecessary blocks that are fetched, are a waste of performance.

Index

S.No.	Title	Page

index in book \Rightarrow find a topic/page faster
index in DB \Rightarrow find blocks of disk faster
 \downarrow blocks that contain
the rows that I
want to fetch.

Purpose of indexes \Rightarrow Reduce # of disk block accesses to fetch data.

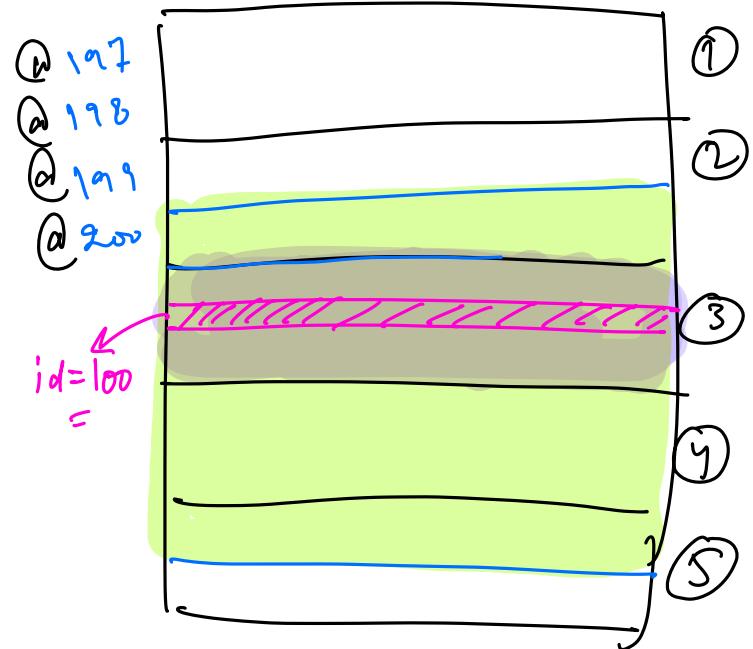
Huge table with 100M rows.

Select * from students where id=100;

Index

id (key)	Block Address (value)
1	1
2	5
3	3
4	4
5	2
.	.
100	3
100M	1

HashMap map.get(100) $\rightarrow O(1)$



w/o hashmap \rightarrow 4 blocks (100M students) \rightarrow
with hashmap \rightarrow 1 block (10000 students)

Select * from students where name = "Akash";

key	value
name	Block address (list)
Akash	[1, 4]
Akhil	[1]
Pallavi	[2, 4]
Rohit	[3]
Arpit	[3]

Students		
id	name	PSP
1	Akash	80
2	Akhil	90
3	Pallavi	78
4	Rishabh	99
5	Rohit	89
6	Arpit	90
7	Vishal	99
8	Akash	98
9	Pallavi	89

Without index \Rightarrow 4 blocks
With index \Rightarrow 2 blocks

HashMap<String, List<Integer>>
map = new HashMap<>();

Goal

- ① You want to clear interviews.
② You want to become a better engineer (Do better at your current job).

find all students with PSP between 60 and 70.

index	psp	block
	90	1
	31	2
	49	1
	67	4
	82	1
	82	4
	90	2
	39	3

Students		
id	name	psp
1	Alok	80
2	Mahi	90
3	Rajani	70
4	Ritikash	99
5	Rohit	89
6	Arpit	90
7	Vijul	99
8	Alok	98
9	Rajani	89

≥ 20 and ≤ 90

for $i=20 \Rightarrow 90$:
if in my map:
`map.get(i)` 20.1
 20.39
 30.41

HM :- check a value $\Rightarrow O(1)$
 get all values in a range $[l, r] \Rightarrow \underline{O(N)}$

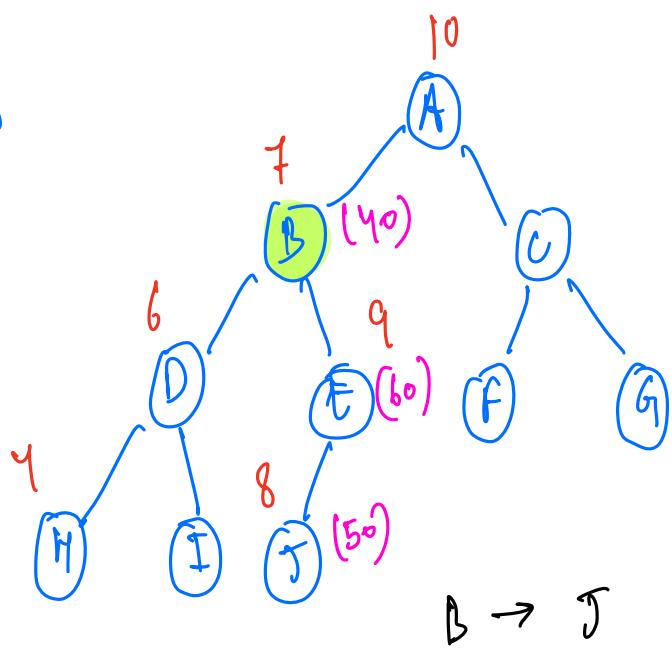
① Sorted

Key
(sorted) \Rightarrow values.

Tree Map \Rightarrow (BBST)
 (Balanced
order-map
Binary Search
Tree)

q.c to check a value $\Rightarrow O(\log(N))$

ques



Given a node,
get the next bigger
node.

$O(\log N)$

B \rightarrow J

40 to 60

① Go to 40

② keep going to next node until you go > 60 .

Internally indexes use \rightarrow B/B+ trees. (optional)

ppp	block
10	1
20	2
30	10
→ 40	9 →
→ 50	3 →
→ 60	6 →
70	70

- ① How indexes work
- ② What is the benefit of indexes
- ③ Why Hashmap is not good.

Break → 10 min

Cons of indexing

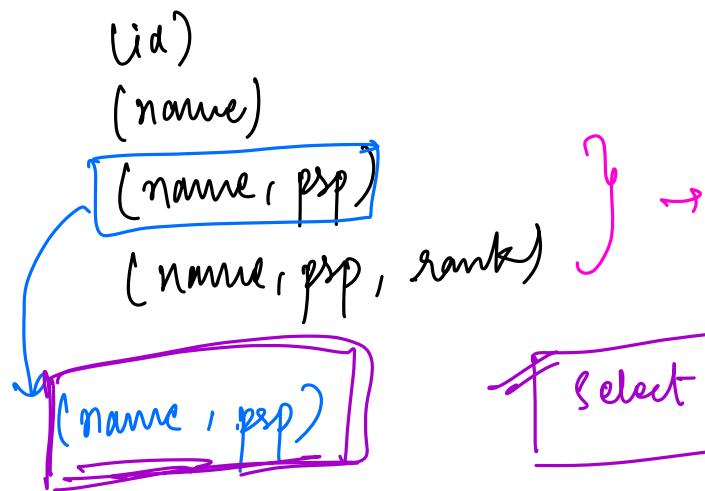


- ① Writes will be slower.
- ② Storage requirements would increase

→ Don't create indexes prematurely

→

Indexing on Multiple Columns.



Select * from Students where psp = 80;

Creating an index
(name, psp)

≠ Create an index on name
and
Create an index on psp

↓
key (Not helpful).

Key	Value	Index Value
(Abhishek - 70)		1
(Abhishek - 90)		2
(Abhishek - 100)		100
(Bipin - 30)		9
(Biswa - 71)		5
(Alok - 20)		6

psp-name (✓)

20 - Alok
30 - Bipin
70 - Abhishek
71 - Biswa
90 - Abhishek
100 - Abhishek

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

Scenarios

query on	index on	help	not help
① name	psp		X
② name	(name) (psp)	✓	X
③ name	(psp, name)		X
④ name	(name, psp)	✓	
⑤ name = X and psp = Y	(psp, name)	✓	
$\boxed{\begin{array}{l} \text{psp} = 1 \\ \text{or} \\ \text{name} = X \end{array}}$		✓	
⑥ name = X $\boxed{\text{or}} \quad \text{psp} = 1$ $\text{psp} = Y \quad \text{or}$ $\text{name} = X$	(psp, name)		

* If a query is on column - X
an index with X as prefix will help.

(X, a, b, c)

(X, a, b)

(X)

$\underline{(a, X)} \rightarrow X$



* $x = a$ and $y = b$
 $y = b$ and $x = a$

$(X) \quad (Y) \rightarrow$
 $\boxed{(X, Y) \rightarrow} \rightarrow \underline{\text{Better}}$
 $\underline{(Y, X) \rightarrow \text{same.}}$

$x = a$ OR $y = b$

(x, y) v/s (x)