# Indexing

## Agenda

Why & what of indexes?

How indexes work

Indexes on multiple columns

cons of Indexes

Indexes on Strings.

select * from film where film-id =100;
↓

TABLE SCAN : row by row ⟶ O(N)

B1    B2    B3   B4
        ↑
      202

HM :        film-id, address.

Index , indexing

O(1) .

select * from users where name = "MOHIT";
                              ↑            ↑

create an
a
index [ key : name , value : row.]

We can have
duplicate name

HM < String, int >  ✗

HM < String, List >

Mohit  |  3 , 7 , 13, 18
Akash  |  1, 10, 16, 18

Select * from film where release-year
between ( 2016, 2023 );

[ Key: release-year, value: List< > ]

2019
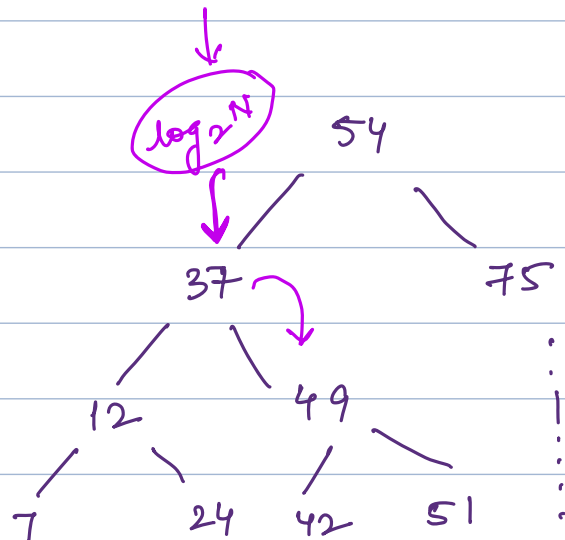2025

range

2016
2017
2018
2019
2020
...

HM's are helpful ⟶ ordered

$log_2 N$

HM + ordered
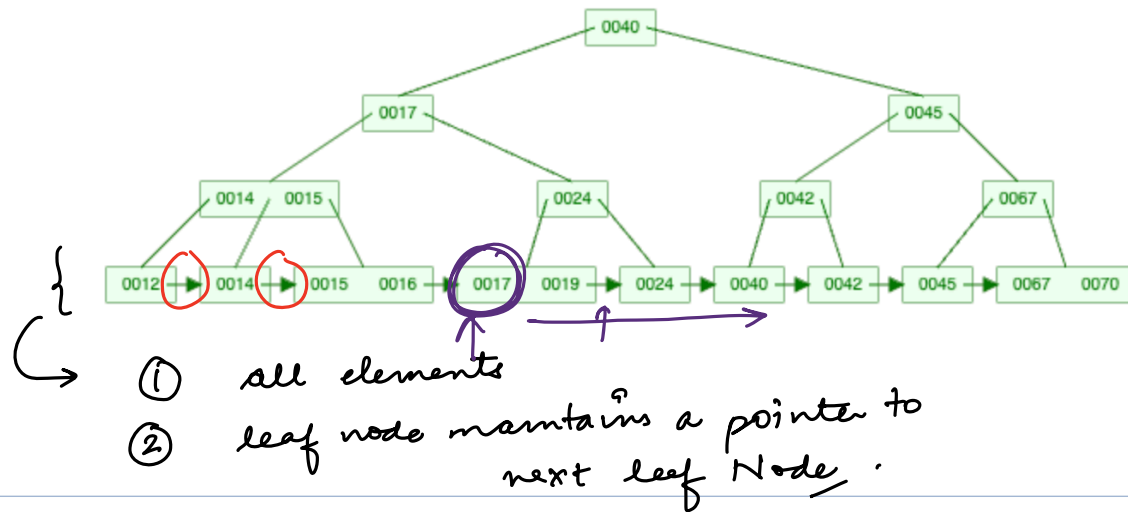
TreeMap/ordered-map
↓
BBST: Balanced Binary
Search Tree .

$O(H)$

$O(log_2 N)$

$log_2 N$   54

37

12        49

7      24    42    51

75

inorder

B / B+ Trees .
↓
decrease height

① all elements

② leaf node maintains a pointer to
next leaf Node.

mysql : tables : indexes : clustered indexes

indexes on primary keys

```
+------+-------+-------+-------+
| id   | a     | b     | c     |
+------+-------+-------+-------+
|   1  | Ag    | B     | C     |
|   2  | Au    | Be    | Co    |
|   3  | Al    | Br    | Cr    |
|   4  | Ar    | Br    | Cd    |
|   5  | Ar    | Br    | C     |
|   6  | Ag    | B     | Co    |
|   7  | At    | Bi    | Ce    |
|   8  | Al    | B     | C     |
|   9  | Al    | B     | Cd    |
|  10  | Ar    | B     | Cd    |
```

| id | a | b | c |
|----|----|----|----|
| 1 | Ag | B | C |
| 2 | Au | Be | Co |
| 3 | Al | Br | Cr |
| 4 | Ar | Br | Cd |

Ag Au Al _Ar_

primary key

| | | |
|---|---|---|
| 1..4 | | Root |
| 1..2 ⟷ 3..4 | | Internal |

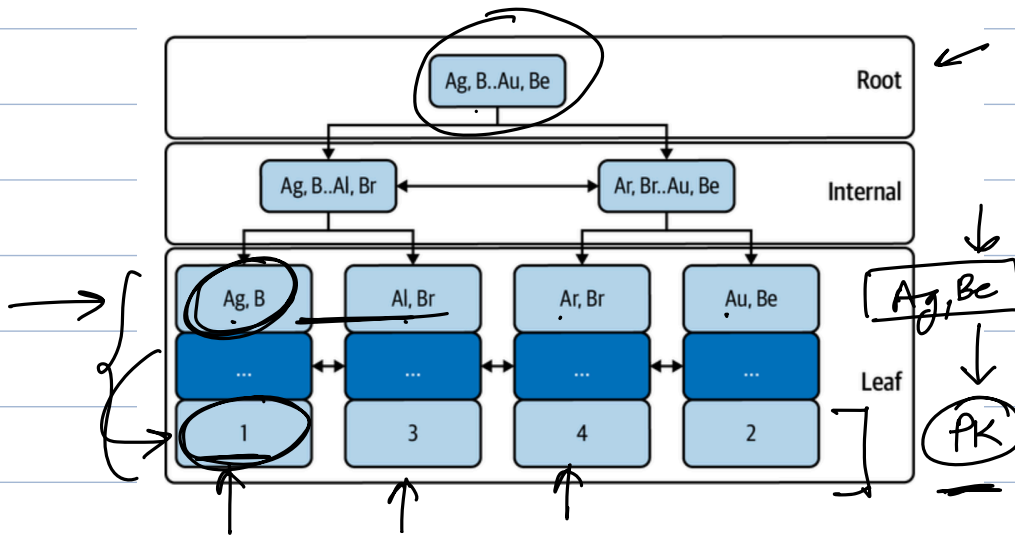| 1 | 2 | 3 | 4 |
|---|---|---|---|
| ... | ... | ... | ... |
| Ag | Au | Al | Ar |
| B | Be | Br | Br |
| C | Co | Cr | Cd |

Leaf

metadata.

select * from film
order by film.id c

We can create our own
indexes : Secondary , Non-clustered index.

select x from Symbol where a = 'Ag' and b = 'Be';

column → PK



Root

Ag, B..Au, Be

Internal

Ag, B..Al, Br ⟷ Ar, Br..Au, Be

Leaf

| Ag, B | Al, Br | Ar, Br | Au, Be |
| --- | --- | --- | --- |
| ... | ... | ... | ... |
| 1 | 3 | 4 | 2 |

Ag, Be

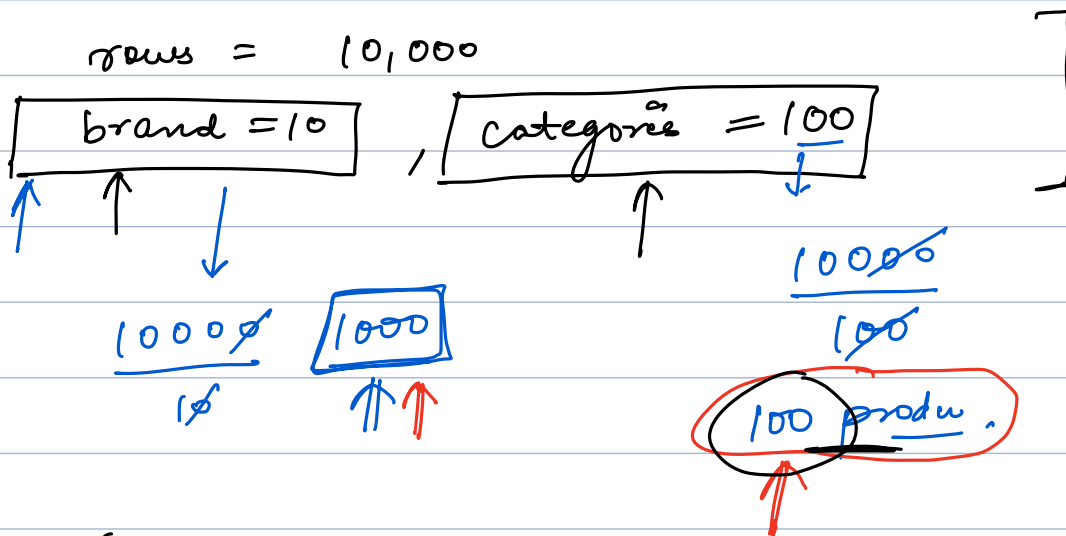PK

Ag, Al, Ar, Au

# Multiple columns

where film-id = ___ ;

order by film-id ;

customers
email, PAN

when email = "___",

& PAN = ___ ;
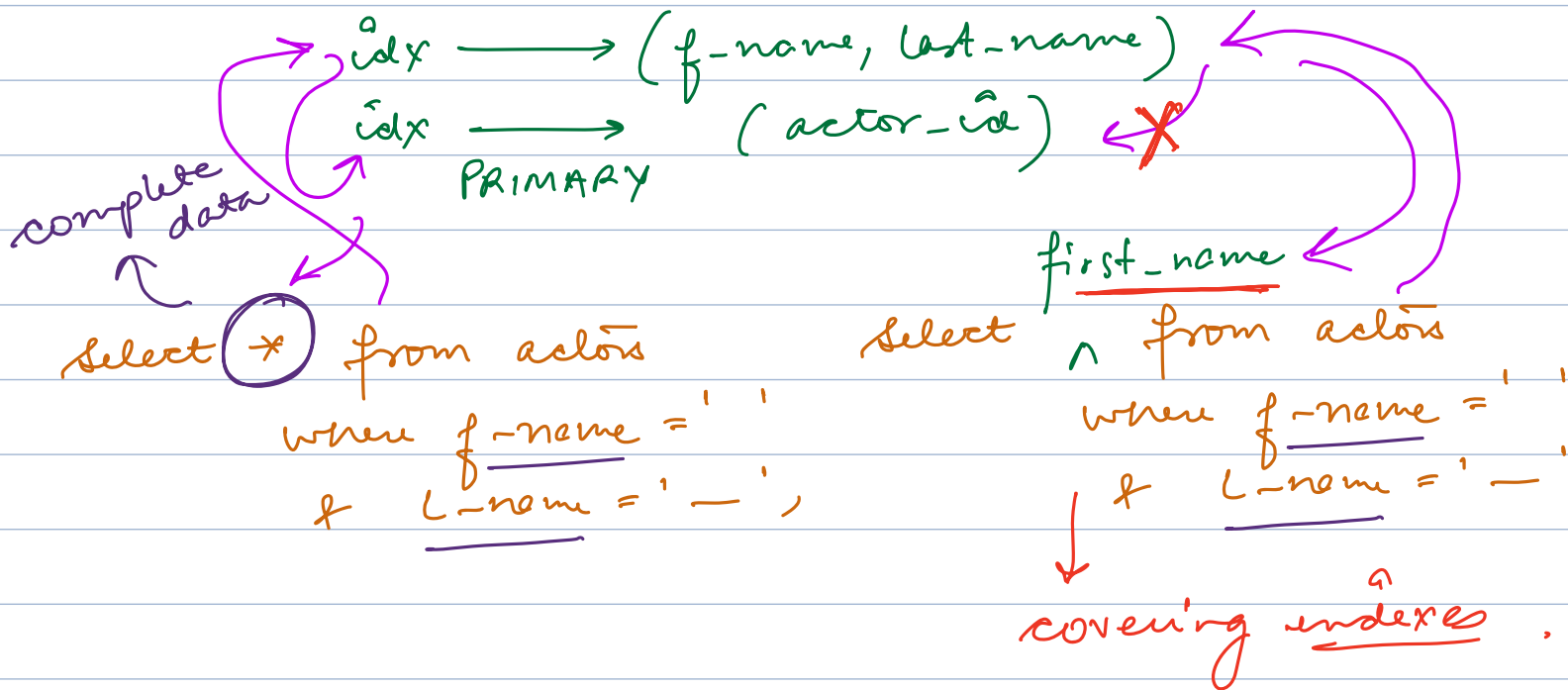
indexes on (PAN) → unique

↓↑

complete table scan.

{ brand
{ category          ——————— brand, category  .

brand, category . X

rows = 10,000

| brand = 10 |  , | categorie = 100 |

$\dfrac{10000}{10}$   1000       $\dfrac{10000}{100}$

$\dfrac{10000}{100}$

100 produ.

(A, B)
                    ] diff indexes
(B, A)              ↓
                    order matters .    { marks, id
                                       { id, marks

idx $\longrightarrow$ (f-name, last-name)

idx $\longrightarrow$ (actor-id) ✗

PRIMARY

complete data

Select ⊛ from actors
  where f-name = ' '
  & L-name = '—' ;

first_name

Select ∧ from actors
  where f-name = ' '
  & L-name = '—'

covering indexes .

I    actor-id

II   actor-id, email

III  last-name, email

LEFT PREFIX

actor-id $\longrightarrow$ I

last-name $\longrightarrow$ III

last-name, actor-id $\longrightarrow$ Ⓘ

email $\longrightarrow$ None of the above .

actor-id, email

| 3 | "a" |
| 1 | "d" |
| 2 | "b" |
| 4 | "c" |

$\longrightarrow$

| 1 | "d" |
| 2 | "b" |
| 3 | "a" |
| 4 | "c" |

$a, b, c, d, e$

I    $a$

II   $a, c$

III  $a, c, b$         c and b    X

IV  $a, b, c$         c          X

a and b

first_N, L_Name          fust-nome

fust-Name

X &rarr; use-less index        last-nome

new idx
will be
required.

where    $a = \underline{\quad}$ and $b = \underline{\quad}$ ;

$a = \underline{\quad}$ OR $b = \underline{\quad}$

$a$

$b \rightarrow$ x Tove

$(a, b)$

$(a, c)$         a     or     b

# CONS OF INDEXING

① indexes occupy space

②       C R U D   ⟶ writes can
           ↓ update     become slower

$a — d$

a mon    bro   Cat   Damon

↓

Momt

12 columns
↓
12 indexes ✓

insert

[BIT MAP INDEX]

low    cardinality :-    no of unique values .
↓

selectivity :    $\dfrac{no\ of\ unique}{total}$ ⟶ $0-1$

drink           1 M row

{ coffee   ⟶   $10^6$ rows .

   tea       2        $\dfrac{2}{1\,M}$

select * from
berey when drink = "tea";

⟶ 0.5 M ⟶ $O(n/2)$
$O(n)$

```
┌─────────────┐
│ Cofee │ Tea │
└─────────────┘
    │         ╲
    ▼          ▼
┌─────────┐   ┌────────┐
│  Cofee  │──▶│  Tea   │
└─────────┘   └────────┘
    │             │
    ▼             ▼
┌─────────┐   ┌────────┐
│ ─ ─ ─ ─ │   │ ─ ─ ─ ─│
└─────────┘   └────────┘
```

( feb 14 )

query ──▶ indexes

update

slow ──▶ locks

93%

# INDEXING ON (STRINGS)

1. bad in comparison
2. space ⇑

_____ @ gmail.com }

where email = "_____";

↑ index ✓

10 M records
↳ 9M

[→] @ gmail.com

create index idx_email on customer( email (7) );

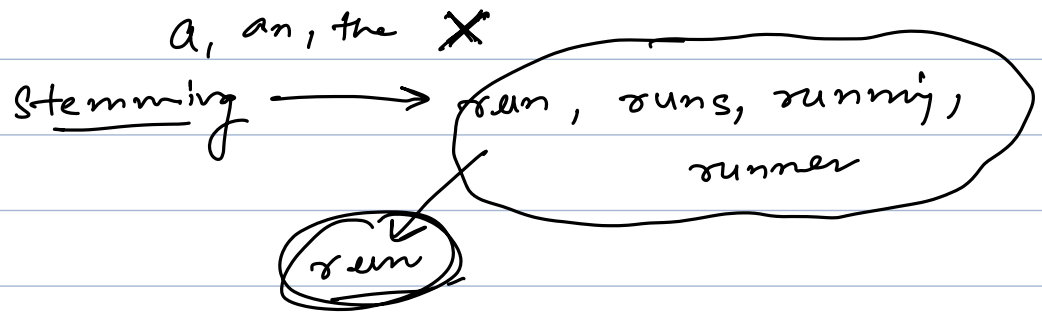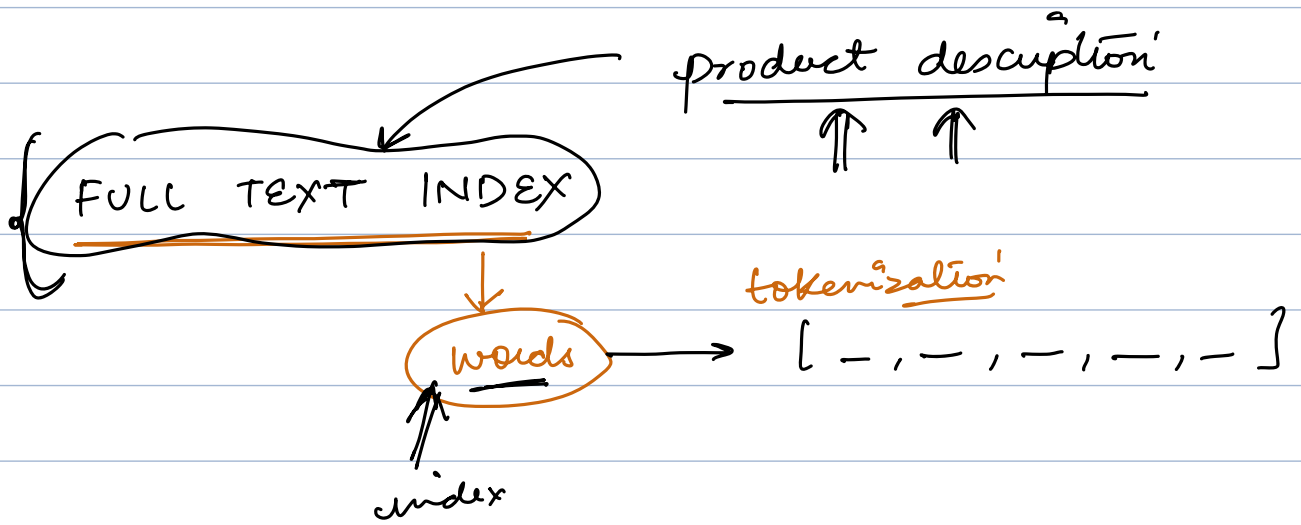2B records

26        26 * 26           17576
         = 676

$26^4 = 4 * 10^5$          $26^7 = 8 * 10^9$
$26^5 = 10^7$
$26^6 = 3 * 10^8$

product description

FULL TEXT INDEX

words → [ —, —, —, —, — ]

tokenization

index

a, an, the ✗

Stemming → run, runs, running, runner

run

MONITOR PERFORMANCE