

NoSQL Databases

~~Not SQL~~

Any way to Store

Data other than

SQL

→ Document

→ Columnar / Column - Family

→ Key - Value Pairs

→ Graph DB

→ Vector DB

All Storing data
in a diff kind

of way with
diff pros and
cons.

→ Relational DB may not be always good

↳ databases other than relational

→ [NoSQL]

Agenda

① Types of NoSQL Databases.

- (i) Key Value Pairs
- (ii) Document DB
- (iii) Columnar DB.

② Case Studies: Situation and you have

to tell what NoSQL Type

will be best for it

Next Class ⇒ How diff types of NoSQL internally work.

① Key Value Pair

→ HashMaps ⇒ (K,V) pairs

K	V
200	"Naresh"
201	"XYZ"
=	

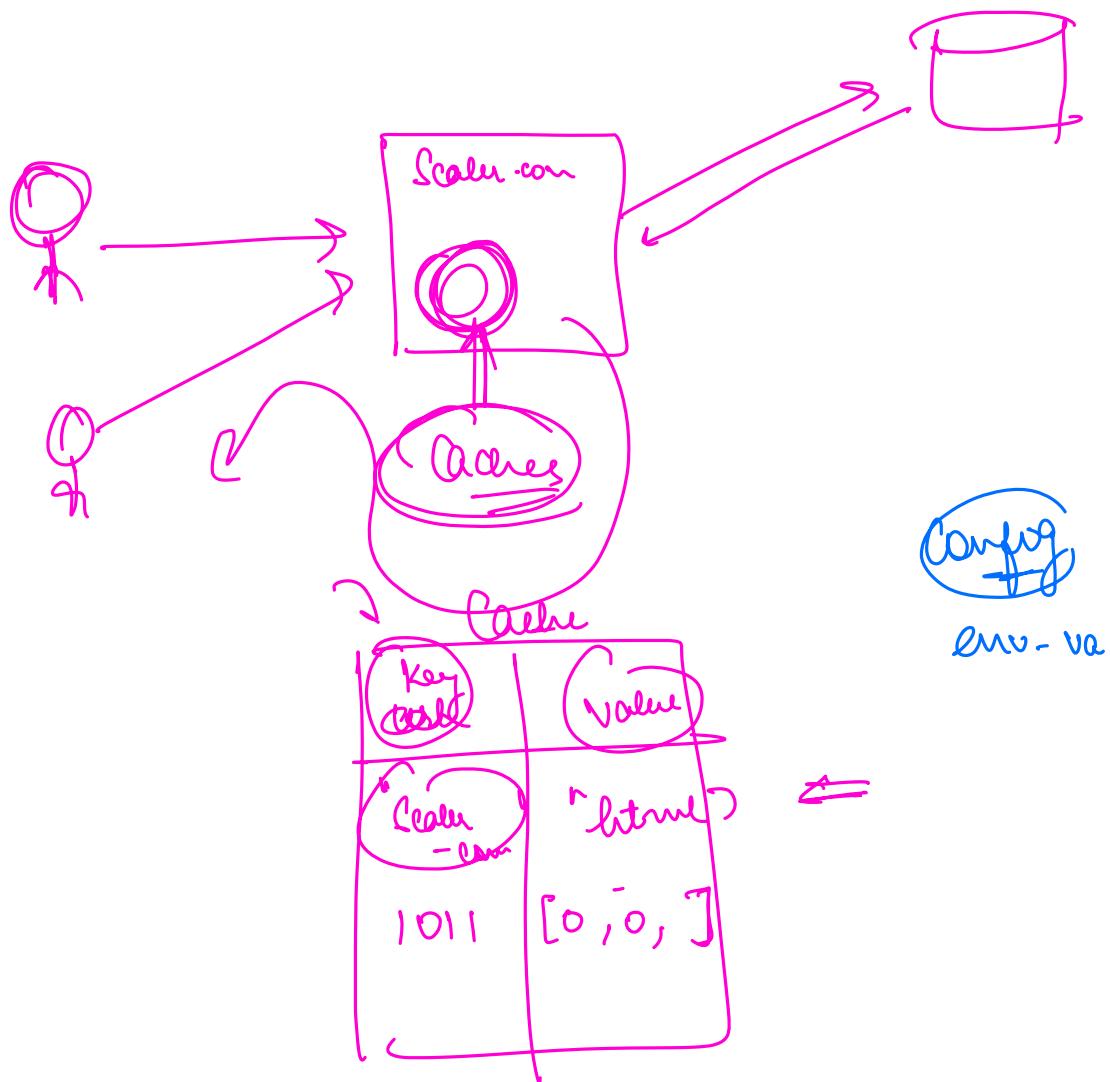
HM[200] ⇒ Naresh

HM[201] ⇒ XYZ

→ HashMaps ⇒ (K,V) pairs
fixed data type

int roomNo	varchar name

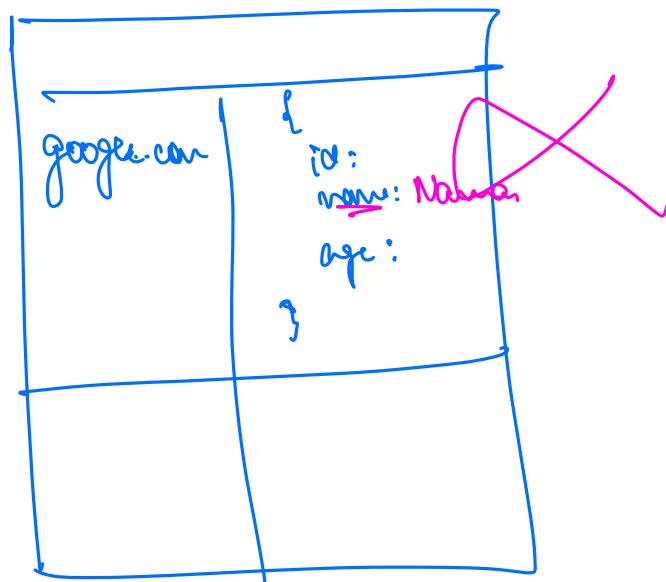
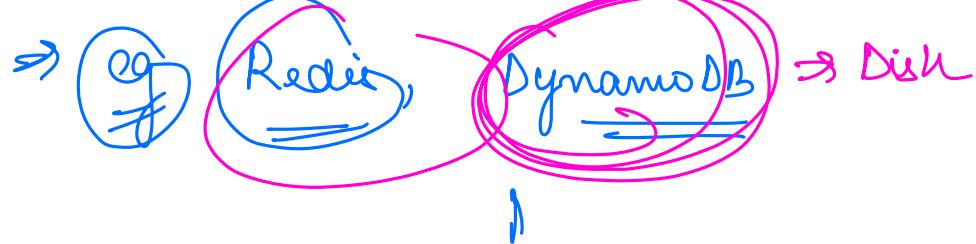
→ typically used when I don't have a fixed data type for values.



⇒ When all the queries are on keys.

= {
 |. get()
 |. set()
}=

work fast in SQL



Type of NoSQL DBs is decided by what type of
it optimizes for.

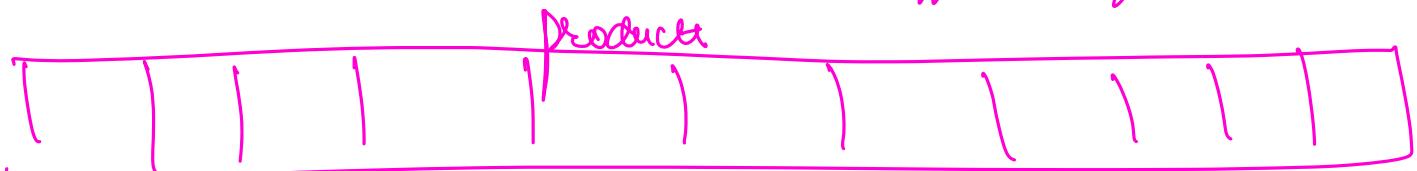
Side Project You can build every type
of DB using a
key value DB

DO Cement DB

0

Issue ⇒ Schema of products
is not well defined

⇒ diff product can have
diff set of attr.



⇒ sparse table

⇒ waste a lot of
Storage

Problem State

→ unstructured data

→ every entity can have its own
set of attr.

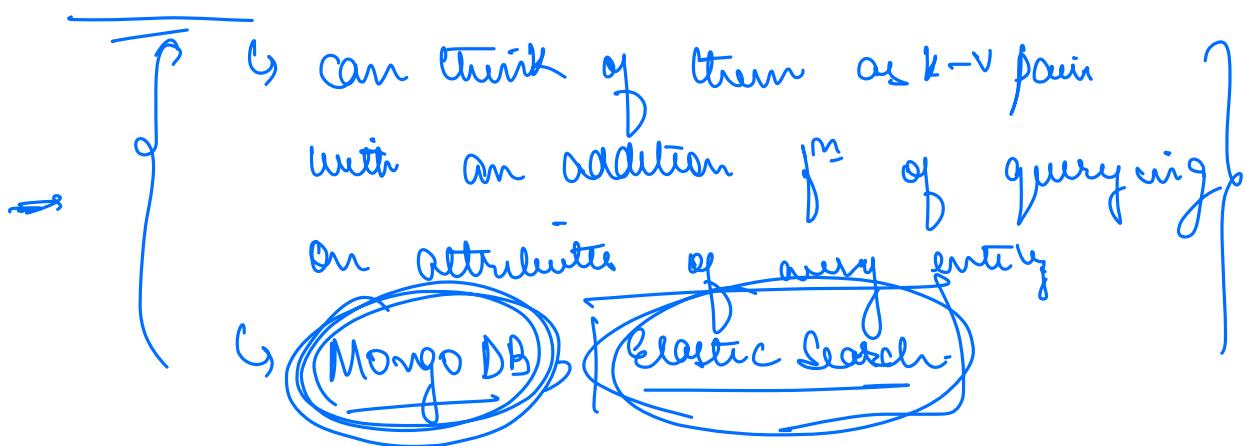
① Ability to store unstructured data

Key / Value

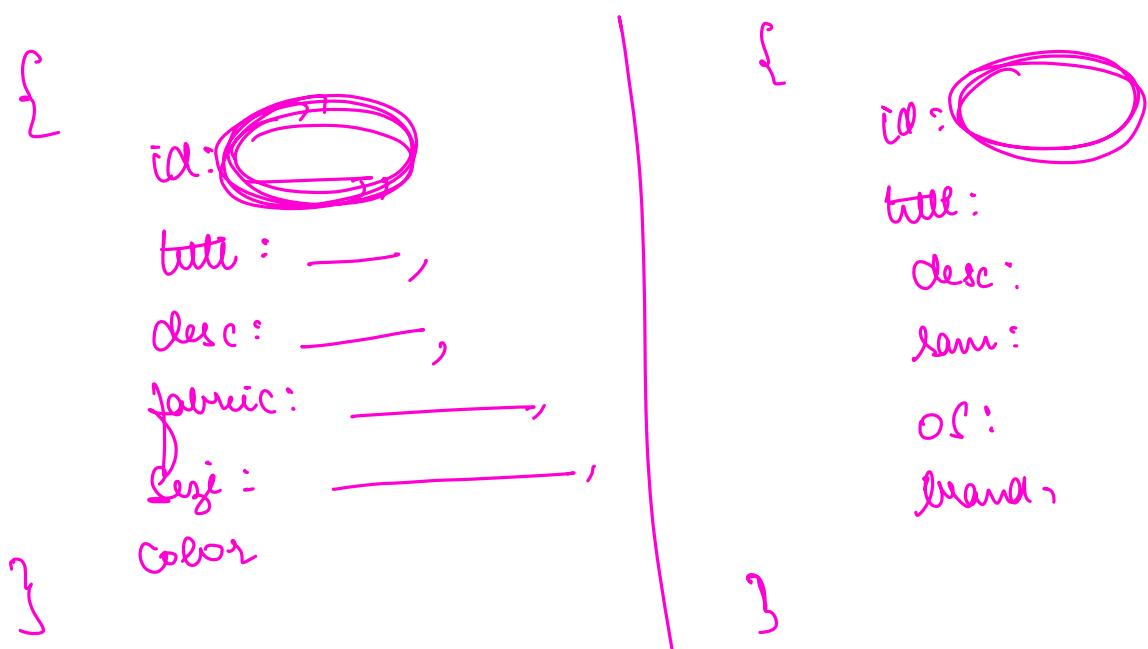
K	V
101	{ id: 100 name: dec type: phone brand: - color: -}

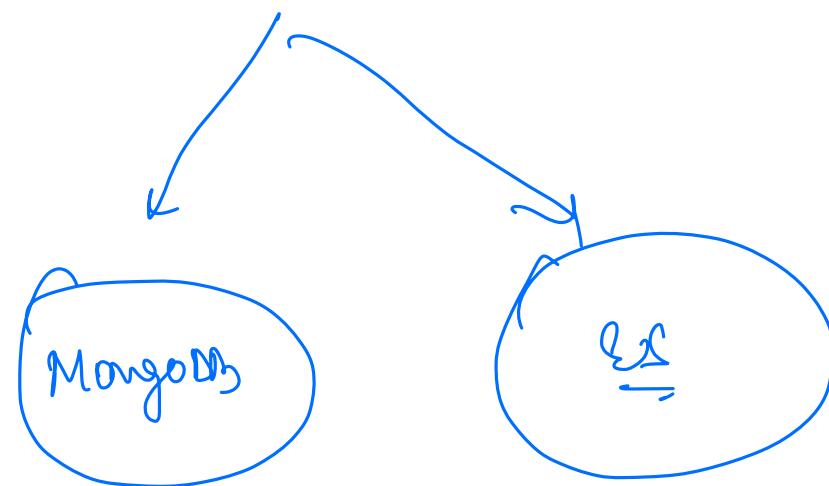
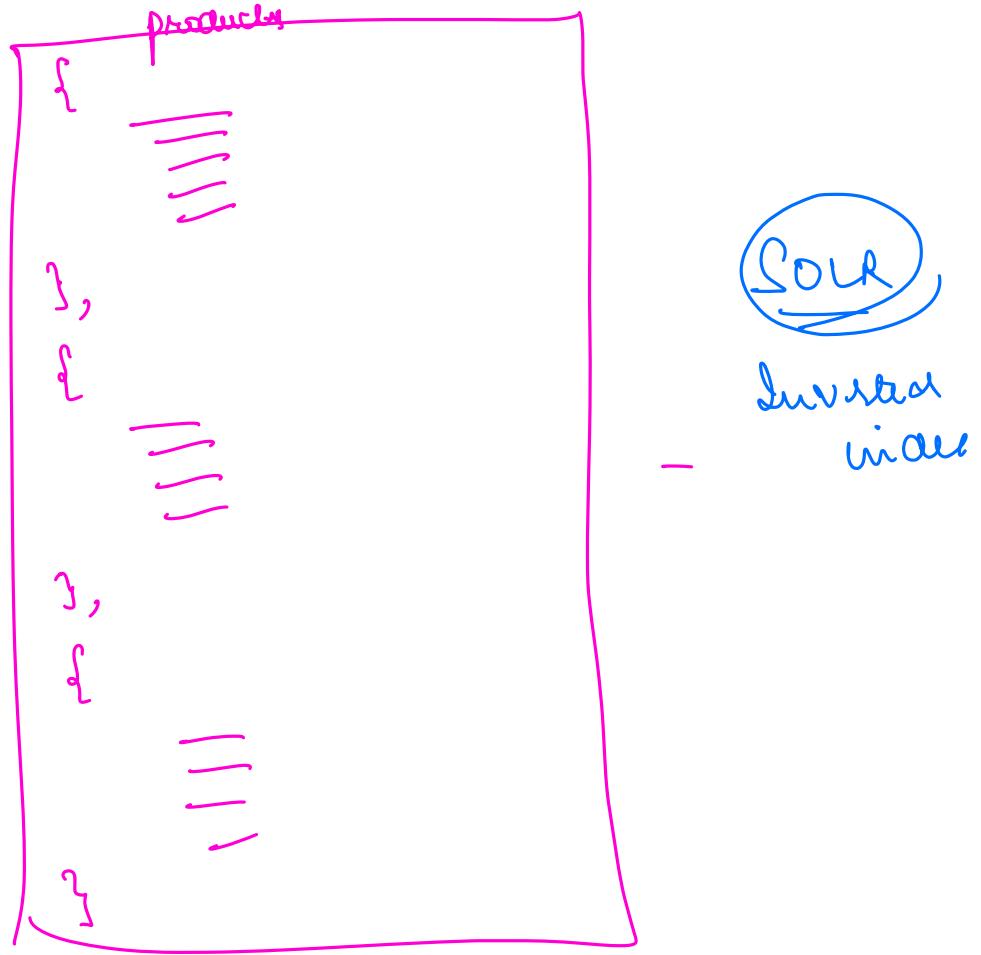
② Ability to query on attributes of unstructured data.

→ Document DB



- ① K-V pairs where values are JSON
- ② every entity is a JSON with whatever attrs it needs to have

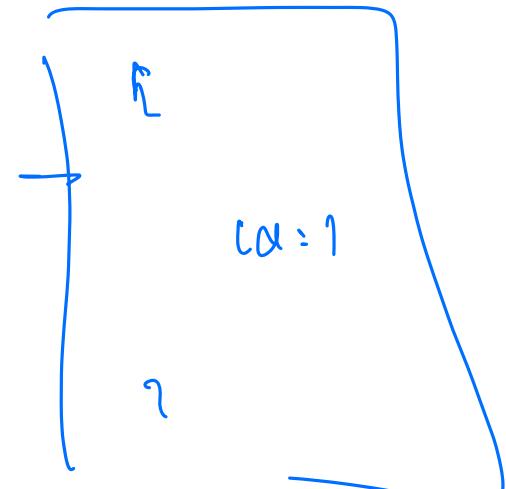




< XML >

< id > 1 < /id >

< name .



Columnar DBs

Column Family

Orders

id	amount	user-id	#of products	date
1	100	a	4	xyz
2	300	b	2	abc
3	200	c	3	def

- ① Select * from orders.
- ② Select * from orders where amount > 100
- ③ Select sum(amount) where user-id = -
- ④ Select sum(amount) -
- ⑤ Select avg(amount) group by date;

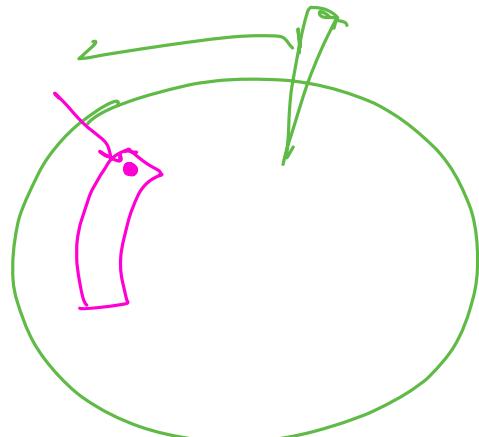
How will SQL internally store data

1	100	a	4	xyz	2	300	b	2	abc	3
	200	a	3	def						

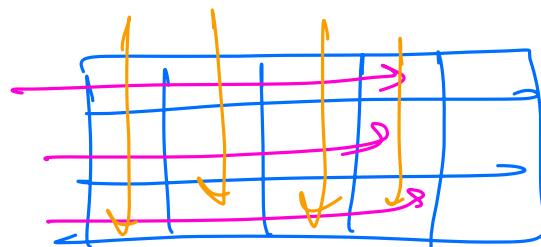
→ index reduce # input/output from disk

disk accesses are slow

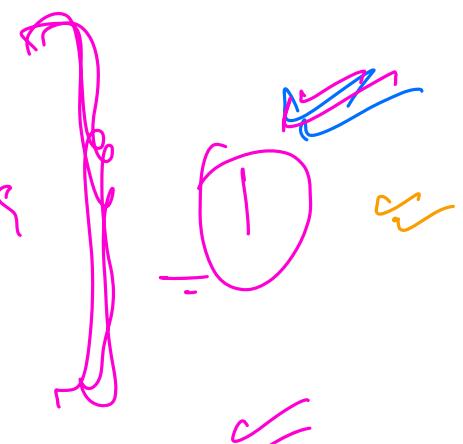
→ when you access a particular data on disk, nearby data is also fetched to mem



a[][]



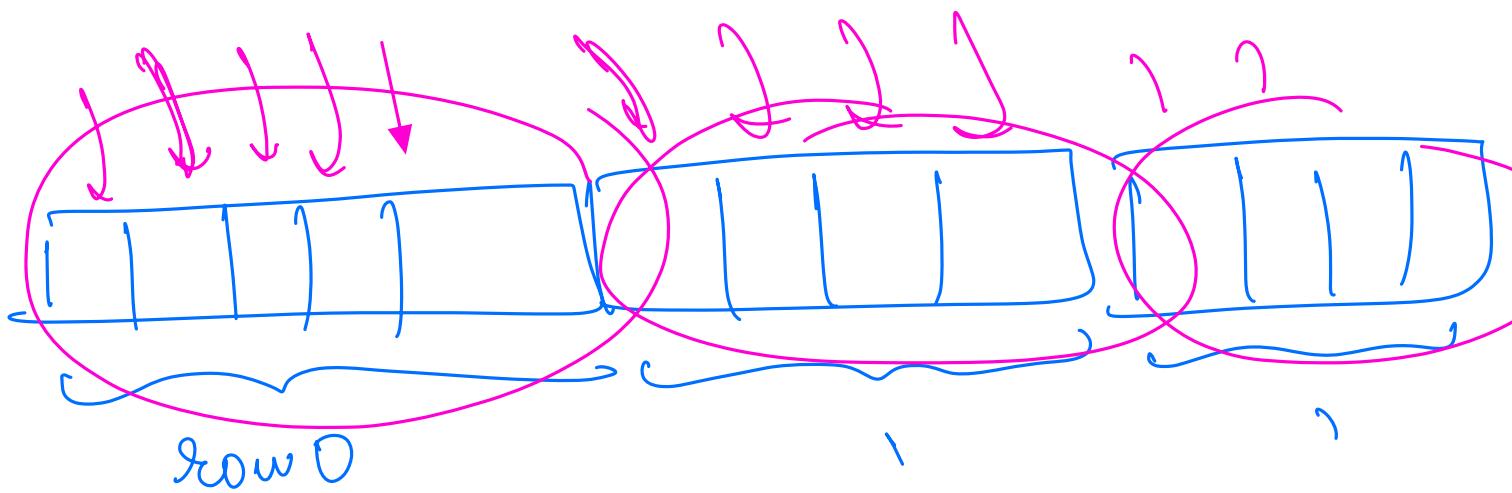
```
for (int i=0; i<n; ++i) {  
    for (int j=0; j<n; ++j) {  
        print(a[i][j])  
    }  
}
```



```
for (int j=0; j<n; ++j) {  
    for (int i=0; i<n; ++i) {  
        print(a[i][j])  
    }  
}
```



② Same



Select k from order

1	100	a	4	XYZ	2	300	a	2	abc	3
2	200	a	3	def						

1	100	a	4	XYZ
2	300	b	2	abc
3	200	a	3	def

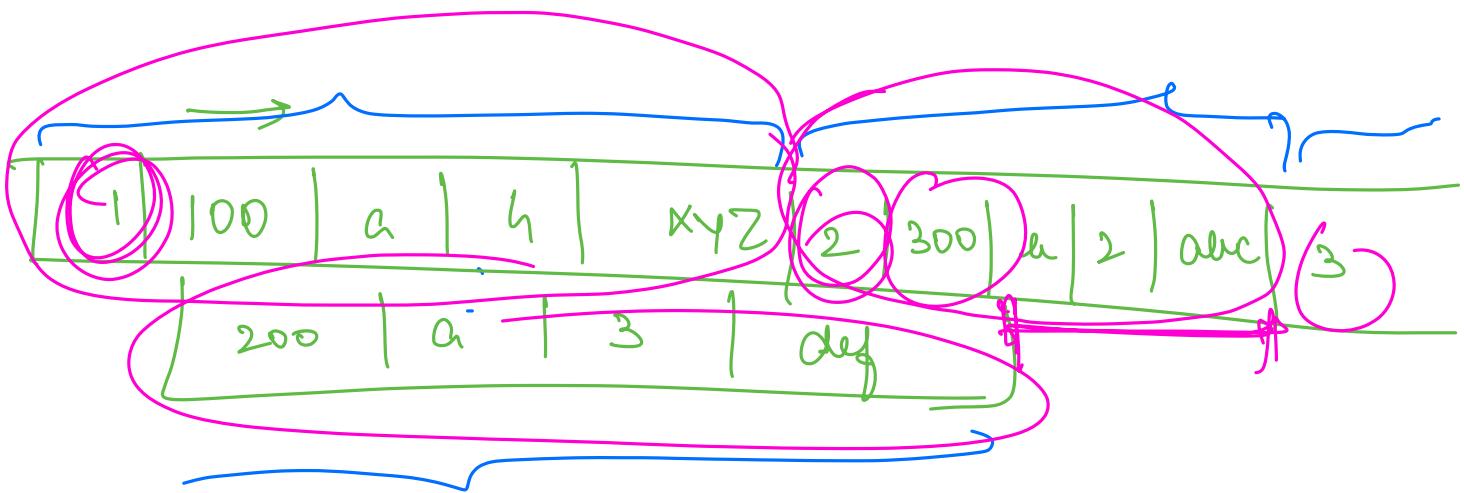
→ All data I ever fetched from disk
was needed.

Select * from Orders where Amount > 100

1	100	a	4	x47	2	300	a	2	abc	3
200	a	3		def						

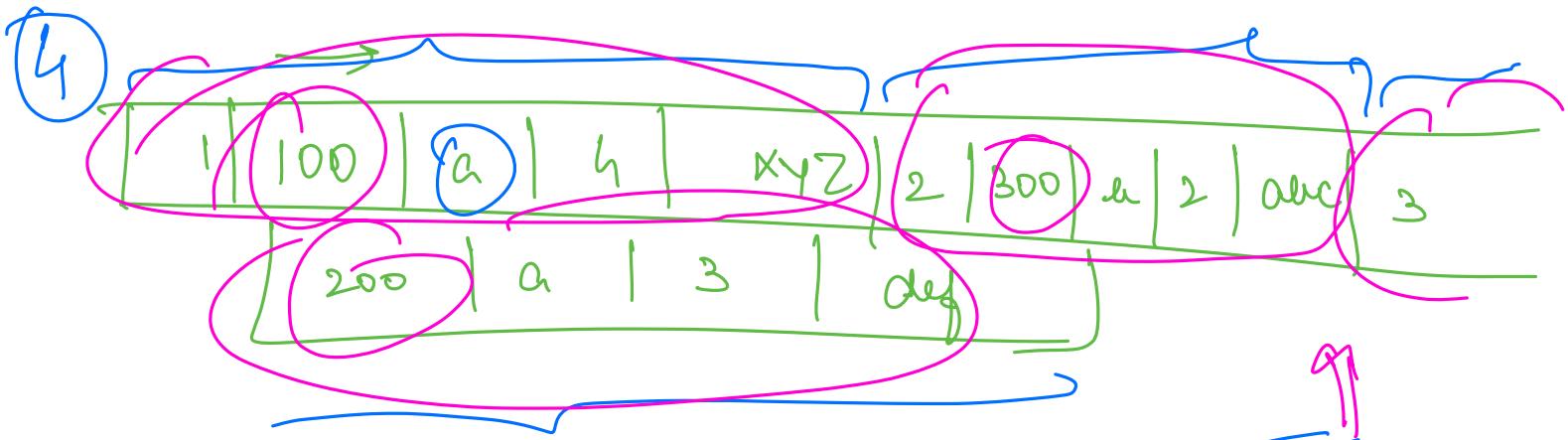
2 | 300 | a | 2 | abc | ~~47~~
3 | 200 | a | 3 | def

⇒ I did some user click access
but fine, it was okay.

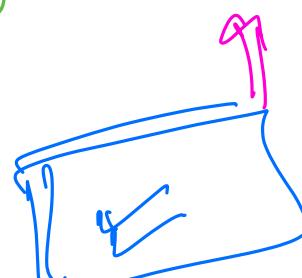


Select sum(amount) where user_id = 2

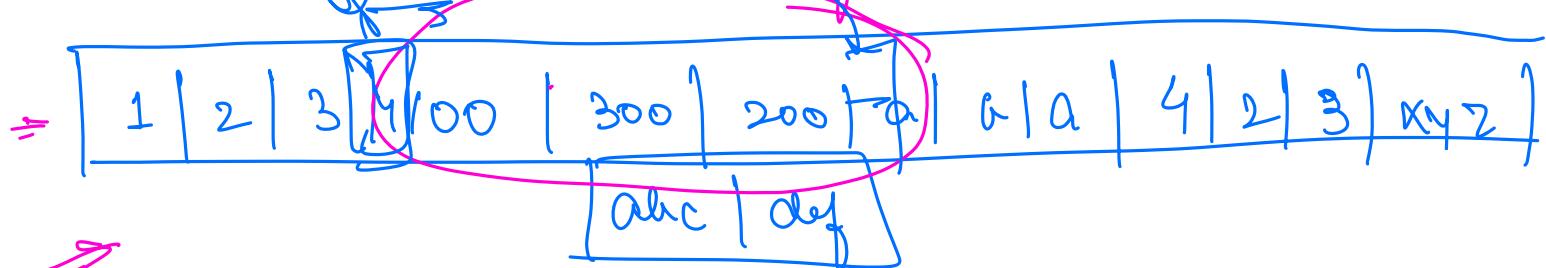
⇒ did I end up fetching a lot of
cols that I never needed.



Select sum (amount)

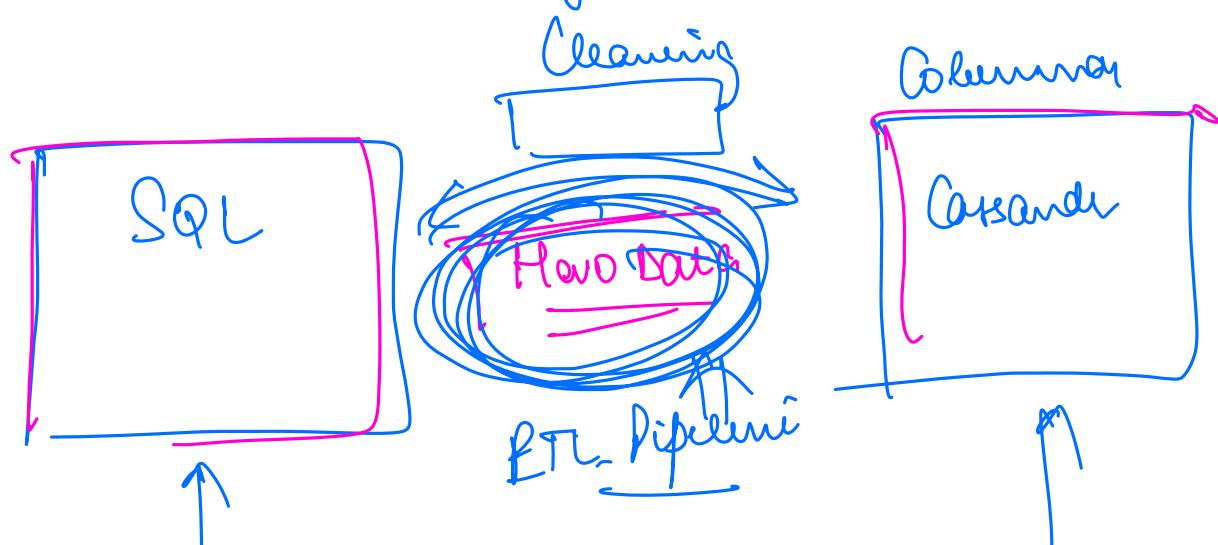


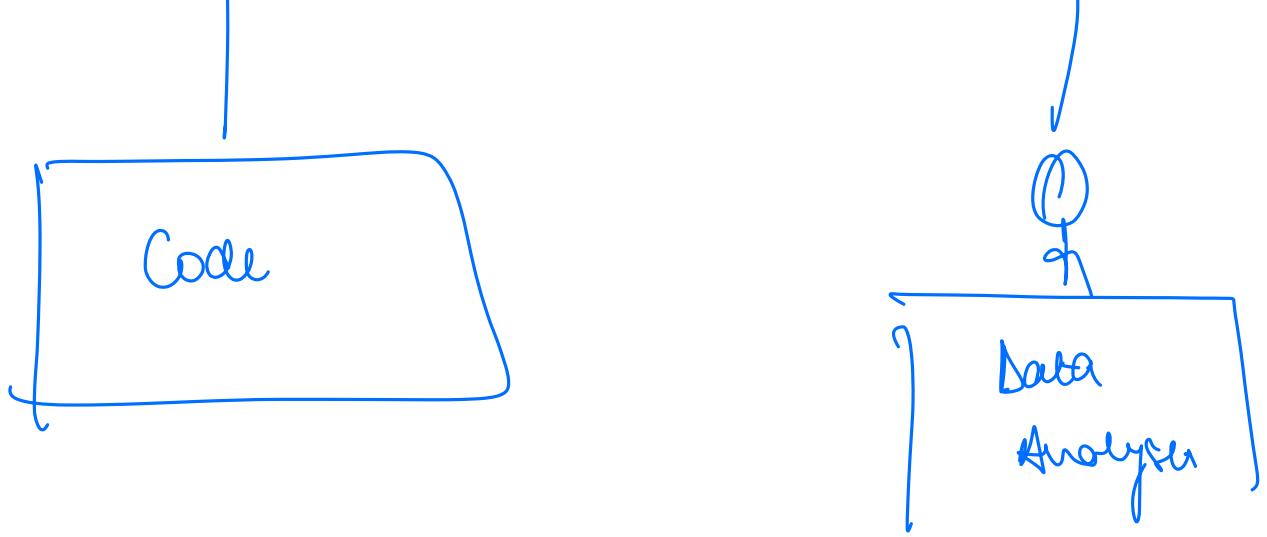
→ I ended up fetching even more unnecessary things



→ Much less disk access!!!

→ data analysis





→ Databases that store data in col^m instead of rows are called columnar DB / column family

Column families

- data is structured
- queries are on value of few columns / agg queries



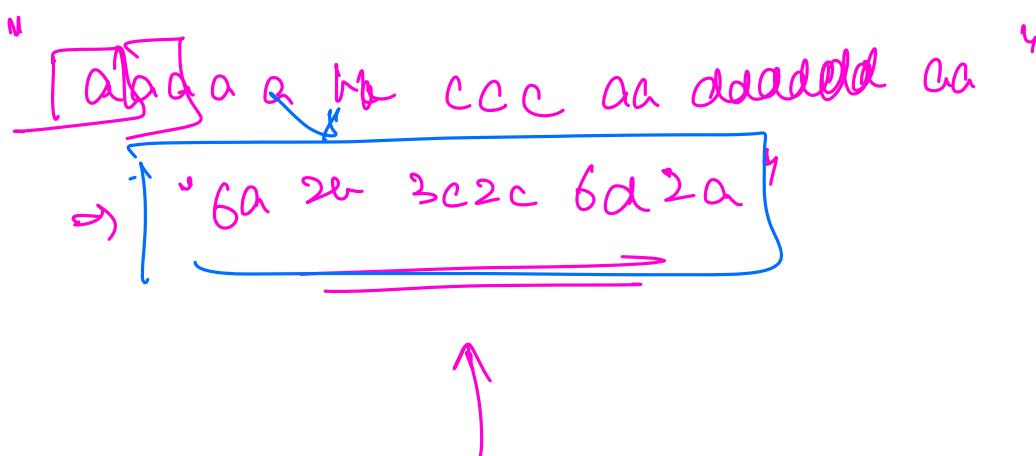
T₁ Cassandra, HBase

Break till 10:30

10:50

- Other optimizations with Col^m family
- Neo4j, Timeseries DB
- Case Study

Qs User



(i) Imagine each character like a DB

row

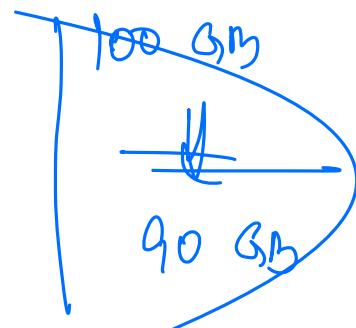
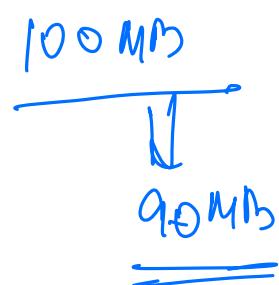
⇒ Can't compact. Because no 2 rows
in the DB are completely same.

products			
id	title	Category	price
1	A	P	100
2	B	P	120
3	C	clothes	140

→ D will store things ordered by cat

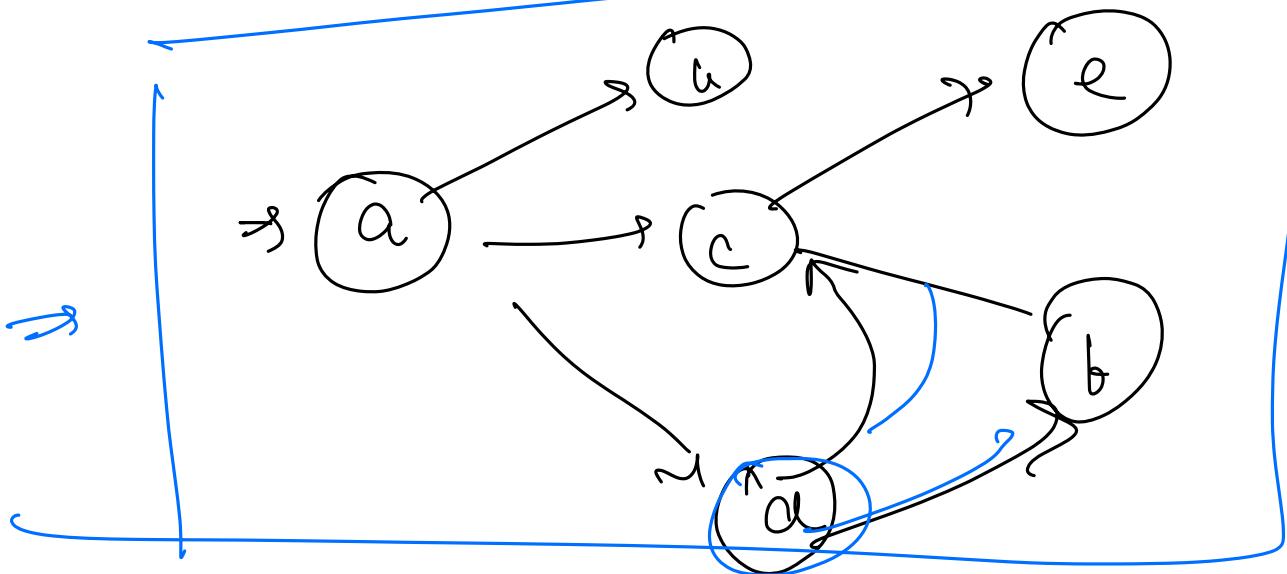
1	2	3	A	B	C	P	P	C	100	120	140	100	100
---	---	---	---	---	---	---	---	---	-----	-----	-----	-----	-----

1	2	3	A	B	C	P	P	C	100	120	140
---	---	---	---	---	---	---	---	---	-----	-----	-----

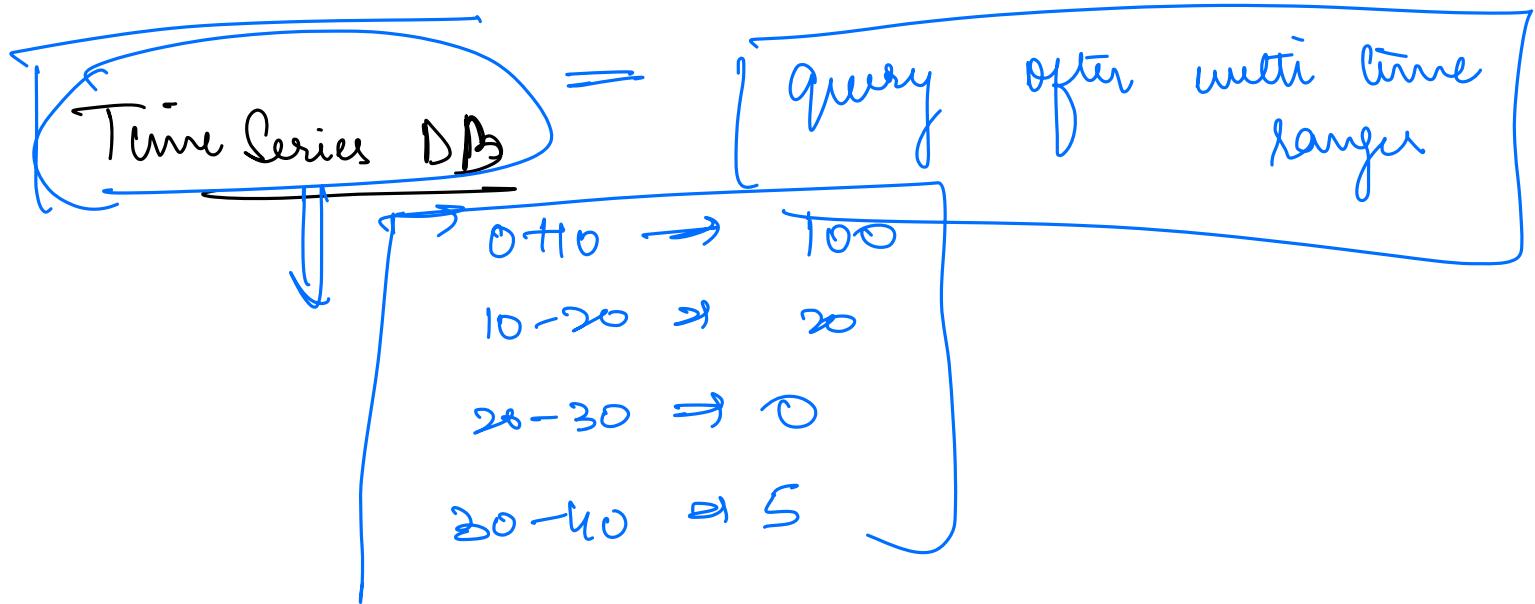


Rel^{^M} b/w entities \Rightarrow Social graph

Graph DB
⇒ Neo4J



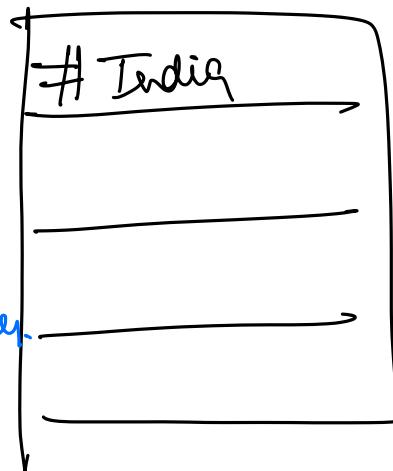
① recommendations



① Twitter Hash Tags

→ a) for each hash tag,
store recent tweets of
that

b) fetching tweets incrementally.



{ 10
↓

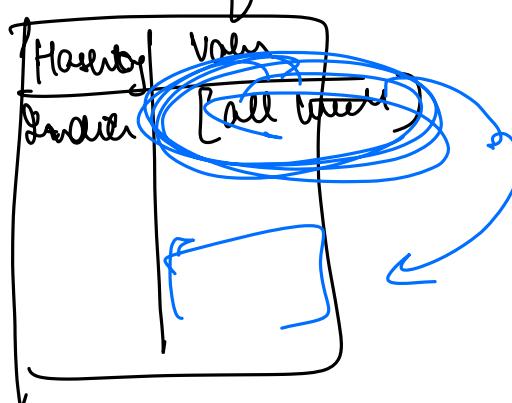
{ 10
↓

{ 10

relational →

key-value →

No. A lot of write queue .



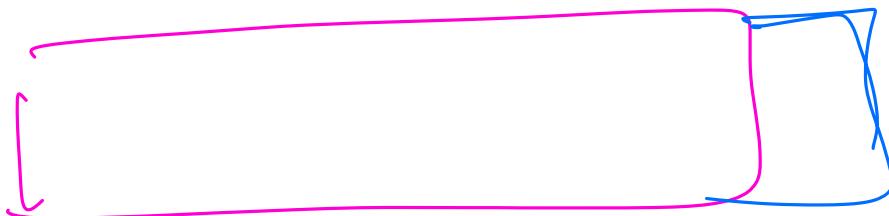
No. Doesn't
allow query
within value

Document → No. Because data is structured

Column-family

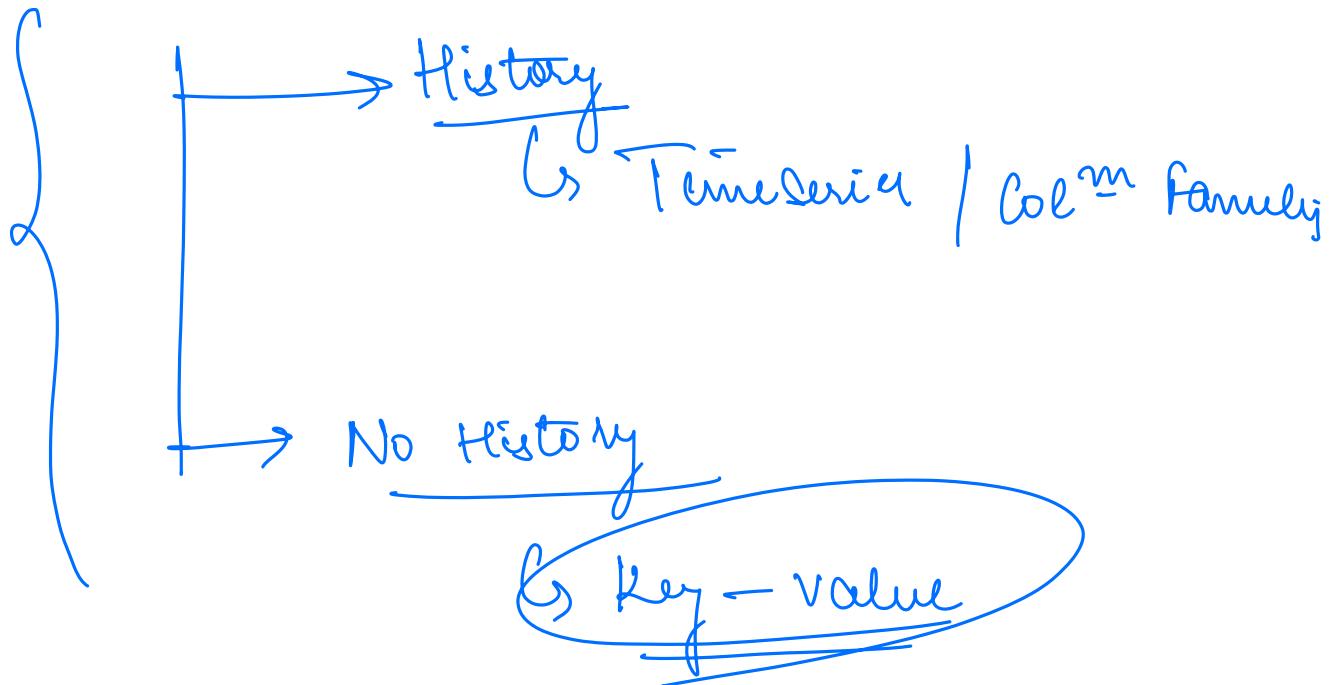
→ ref

3. If writes are append only they can handle very fast.



Live Score of a Match (Current live score)

$$\textcircled{101} = \textcircled{121}/8$$



Current loc^m of an Uber Driver

- ↳ Key - Value (if history not reqd)
- ↳ Columnar / Time Series (if History)
- ↳ Geo spatial DB

⇒ ↳ get me drivers in 2 km of [lat, long]

↳ get 7 nearest driver