

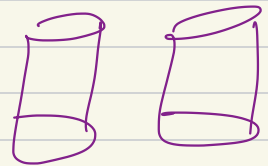
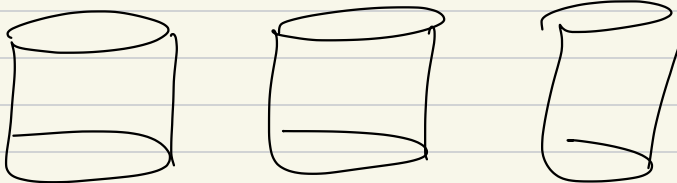
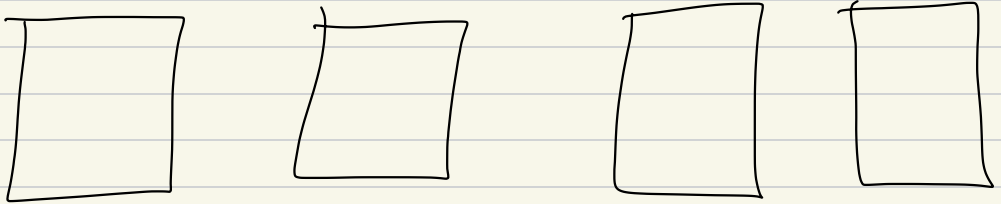
22 / Dec / 2023

STORING LARGE FILES

Quad Trees

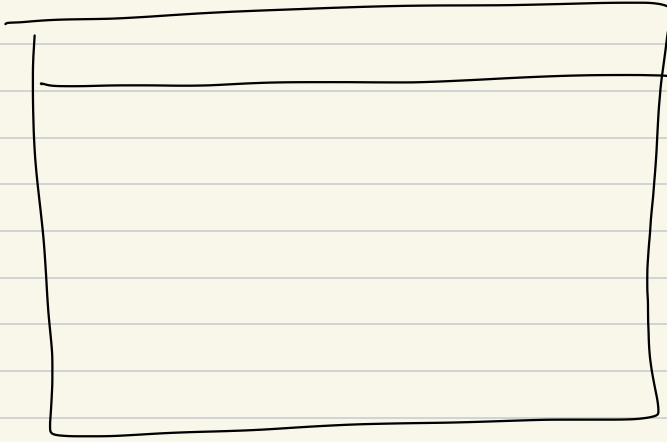
Client

Gateway + Load Balancer



File Store / Blob / Object

Why NOT a .SQL DB?



NoSQL X

distributed file system 😊

HDFS

AWS S3

Ep

Scaler

3GB X 10,000 Videos

$3GB \times 10^4$

30TB

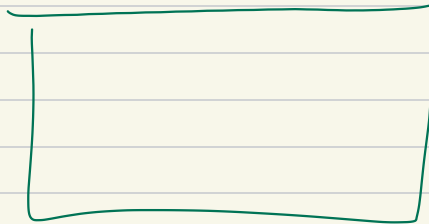
→ replicate

HDFS



or

AWS S3



ex.

few GB

~5GB

x

100k

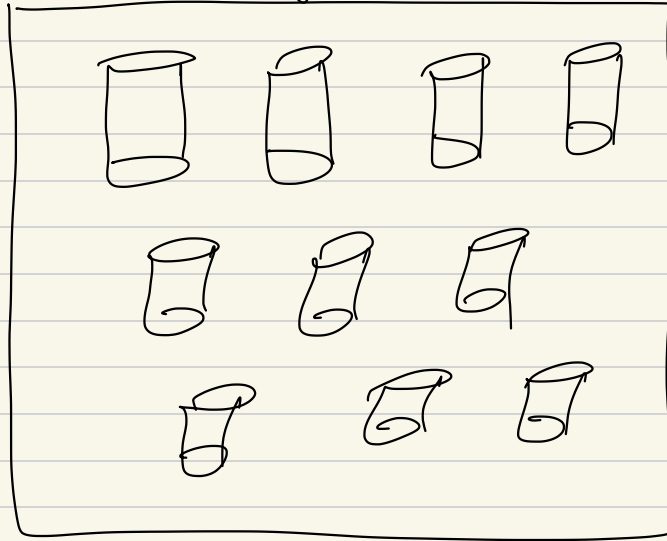
≈

100TB

Files that are big

- Text Files
- Log Files
- media files
 - videos
 - images
- neural networks

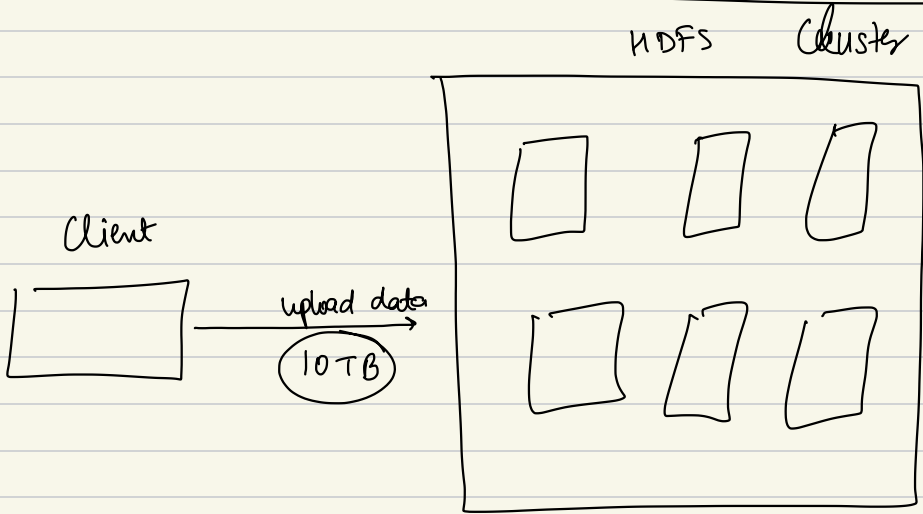
cluster of machines



- ① Huge data storage capability
- ② Data Availability
- ③ Data Uploading should not require us to restart from 0 once network fails temporarily.

④ Downloading

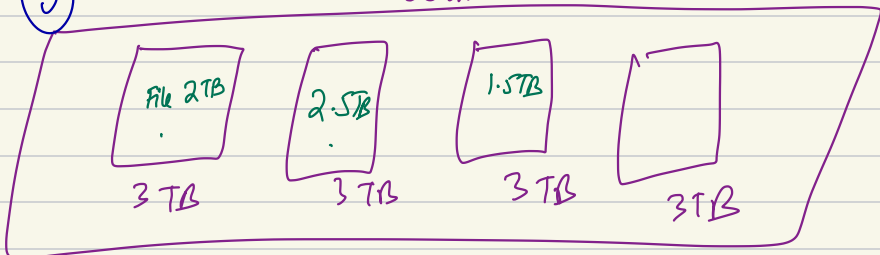
— not require to restart
from 0 after n/w failure



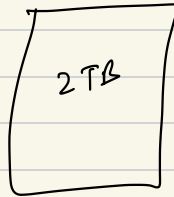
Problems

① 😊 1 File itself might be too big to save
in 1 machine.

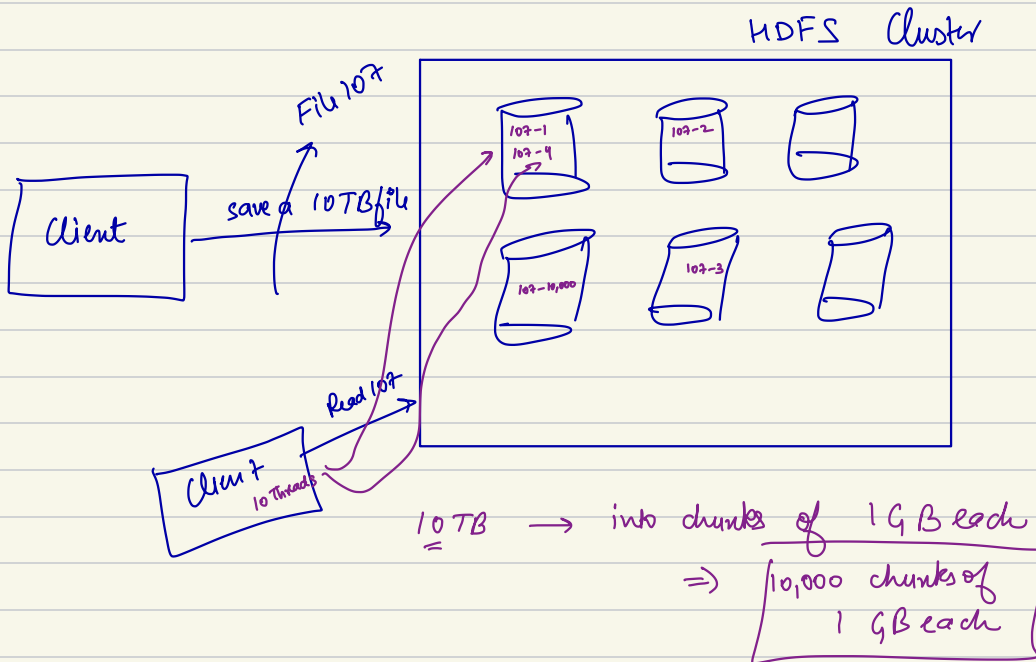
② 😊 Fragmentation } under-utilizing your
cluster space



③ Parallel downloading is not possible
if you store entire file as 1 unit



Files are going to be chunkified



10TB



chunk size spectrum

Tradeoff

↑
simply too many chunks to manage

$$\frac{10^9}{10^{12} \text{ file}} = 10^3$$

10^{10} chunk

0,00,00,00,000

→ 10 core chunks

metadata ☺

File Id

107
107

chunk Id

1
2

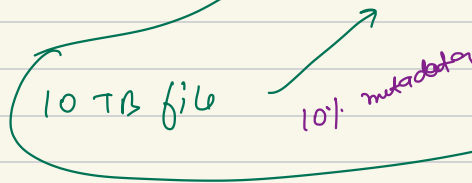
Location

Machine B - path ⇒ 100B
Machine X - path

10^{10} rows file

$$10^{10} \times 100 \text{ B}$$

$$= 10^{12} \text{ B} = 1 \text{ TB of metadata}$$



HDFS

default chunk size

1.0	→	64 MB chunk
2.0	→	128 MB chunk

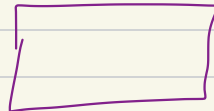
hw:

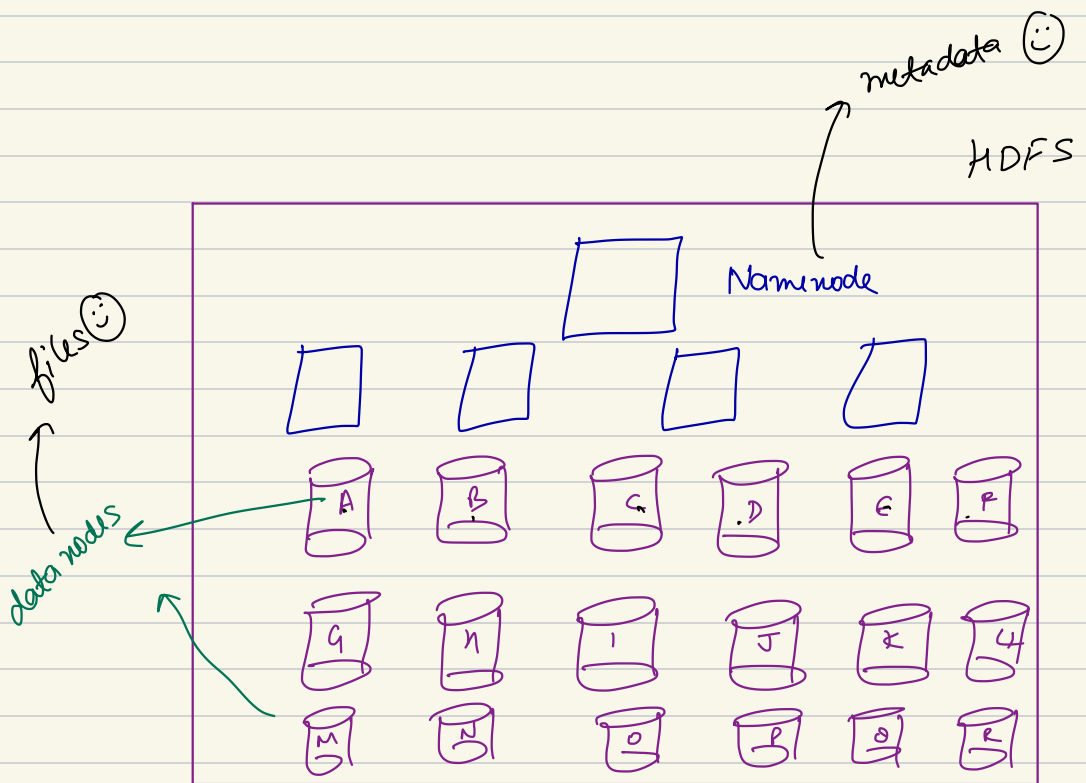
10 TB divide 128 MB chunk



← metadata 100 bytes

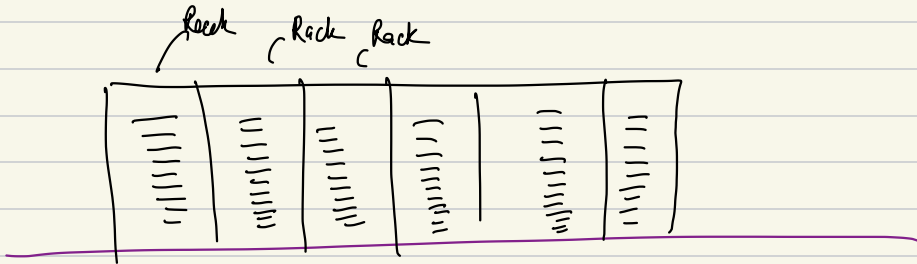
★ metadata = ?



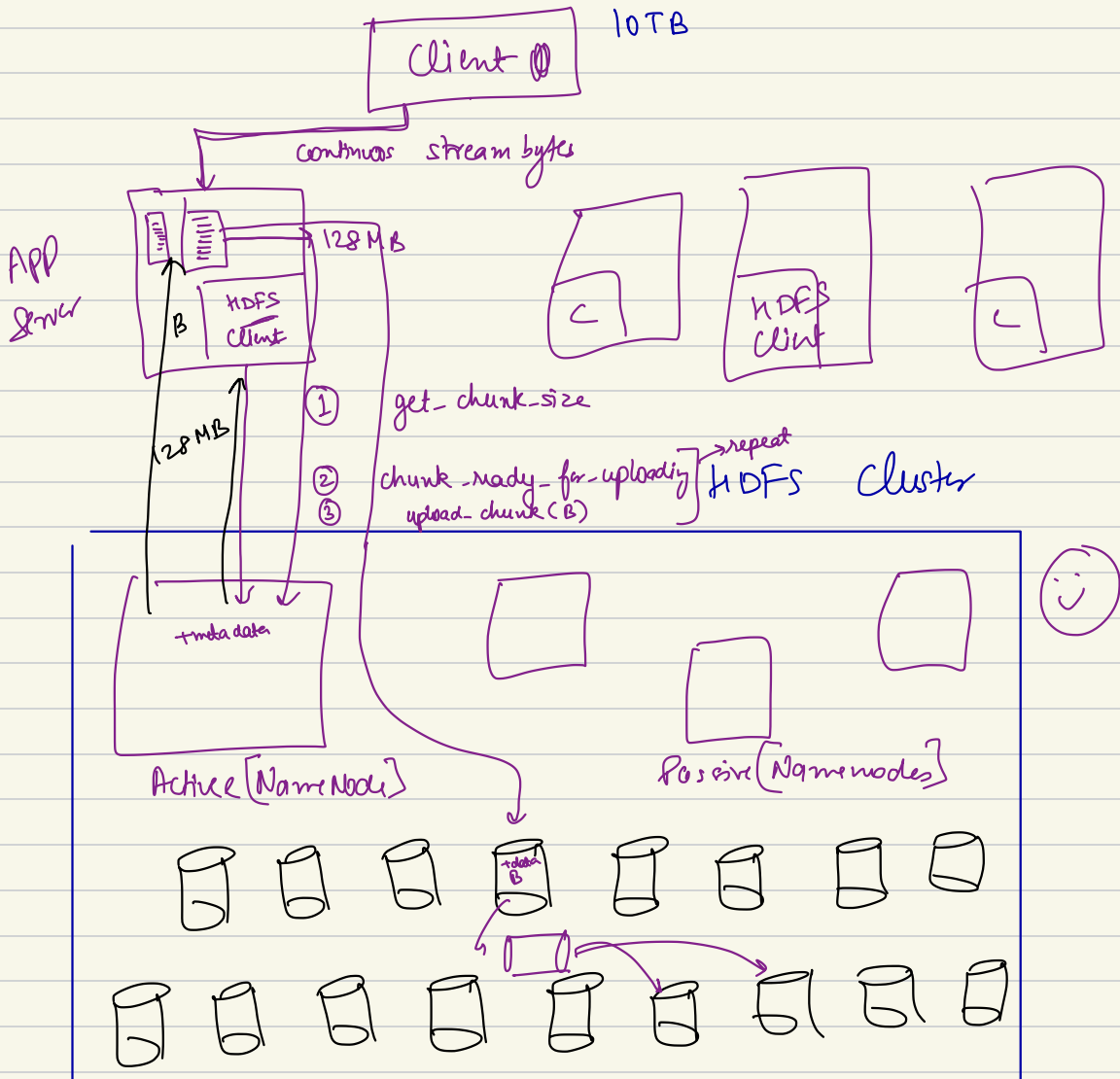


Most Reliable 😊

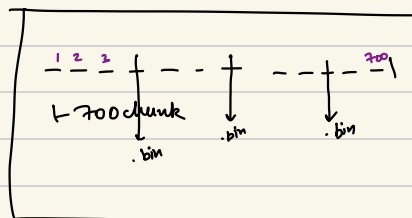
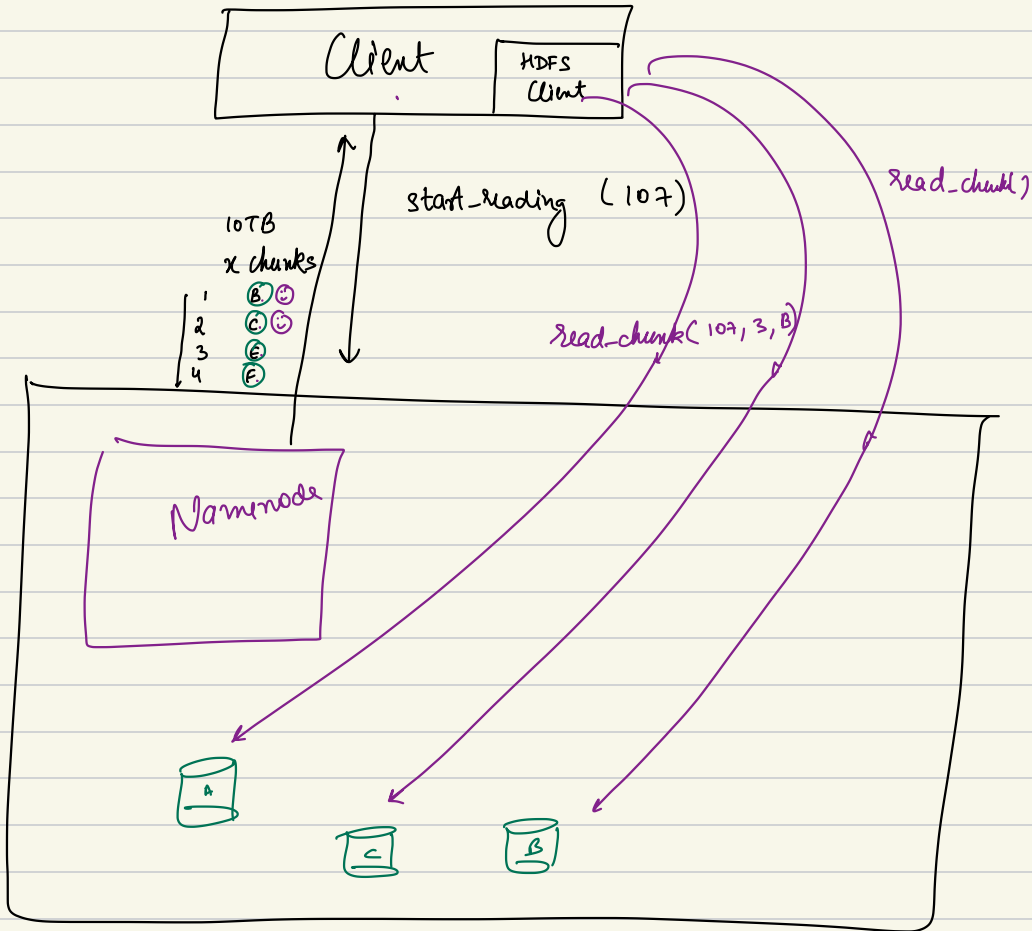
Rack Aware Algorithm 😊



Write Flow | Upload Fl

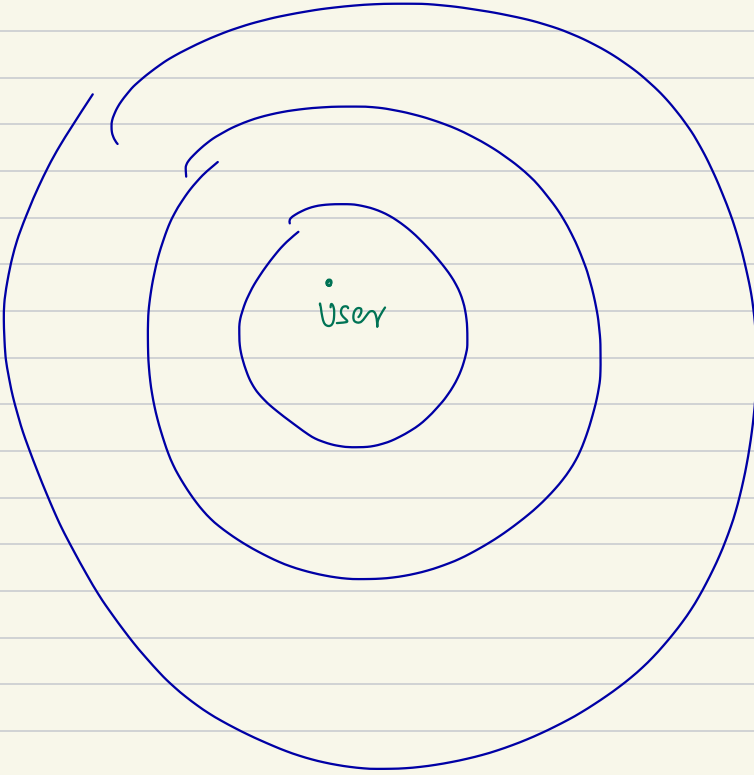


Read Flow | Down load



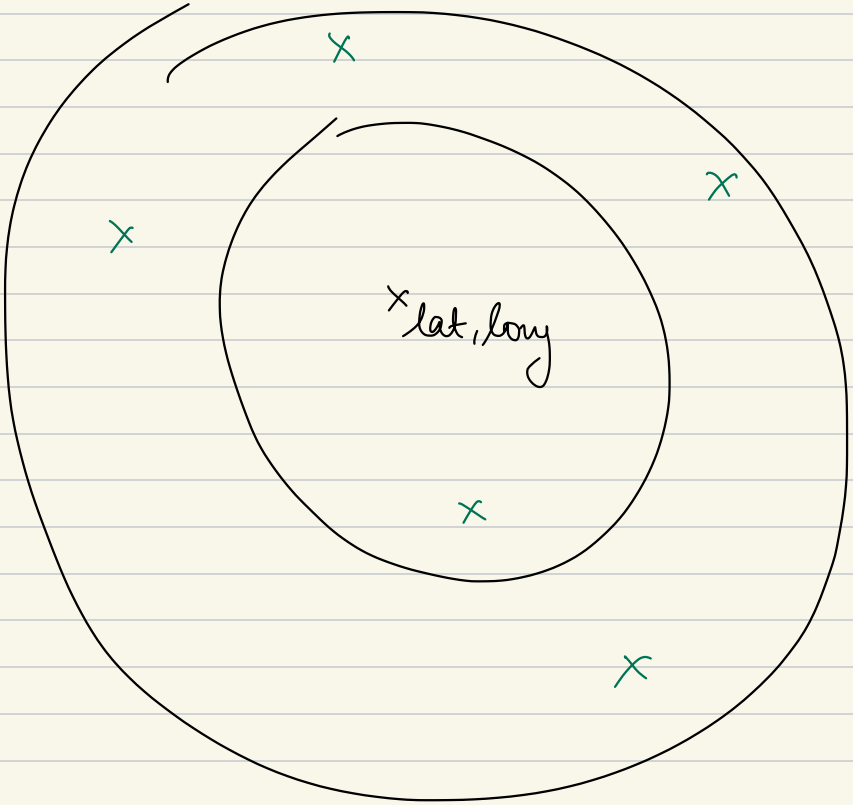
VAD TREES

Uber



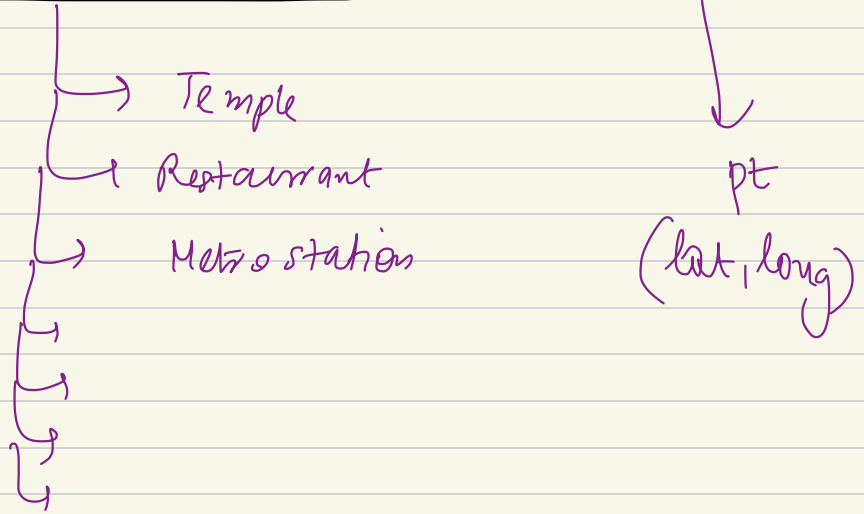
google maps

barber shops near me



Places near me

Place of interest near me



	Id	Type	Name	Location
Step 1 :	107	<u>Barbershop</u>	Alim Hakim Barber	(lat, long)

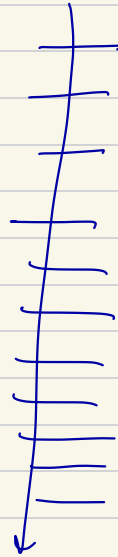
100M Places of Interest

Query :

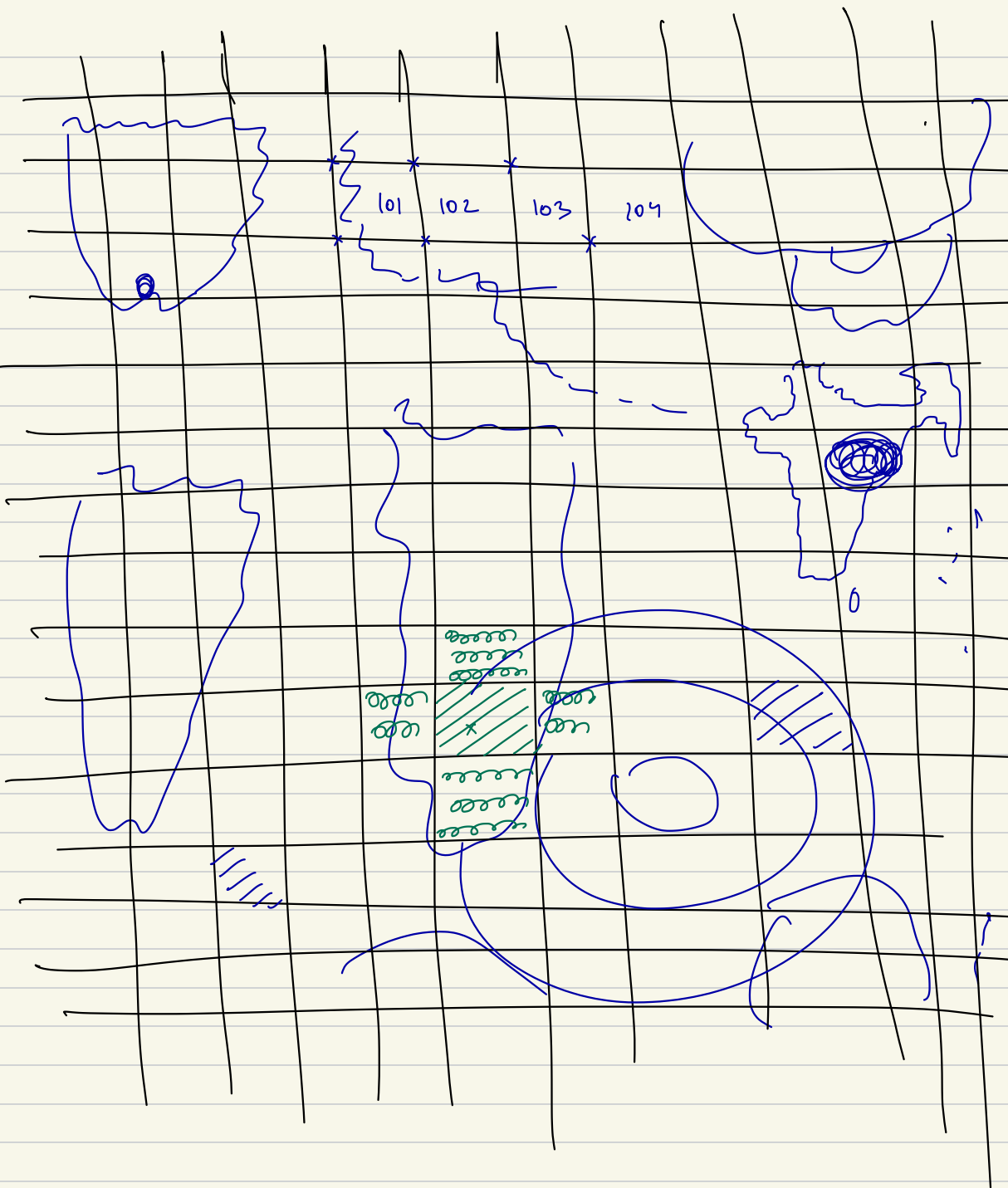
X
[lat, long]

nearby petrol
pumps

U_x, U_y



$$\sqrt{(U_x - L_x)^2 + (U_y - L_y)^2}$$



Iteration
1

Simple indexing OR sharding will not solve the issue as; even after it you will have to consider 10s of Thousands of comparisons.

⇒ Normal SQL / NoSQL DBs won't work 😞

Iteration
2

Divide the world into grids; and then based on the query, identify grid-id and then perform a place of interest search only in that grid cell
[+ neighbouring grid cells]

PROBLEM:

Even though a grid cell solution is a solution in the right direction, but it is NOT without problems

Unequal density of places of interest
throughout the world

Flexible Grid Solution

