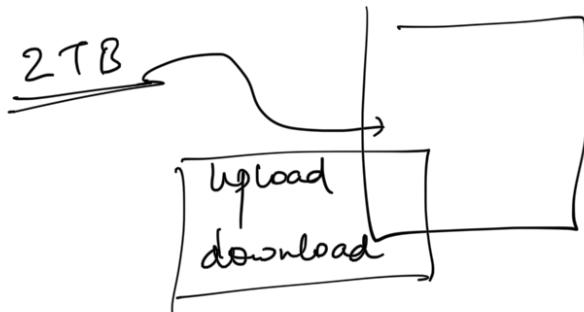


S3 QUAD TREES

HDFS

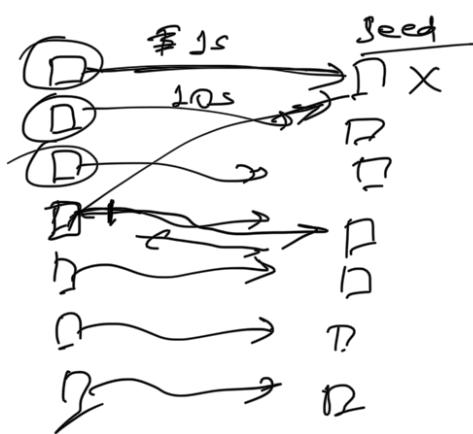
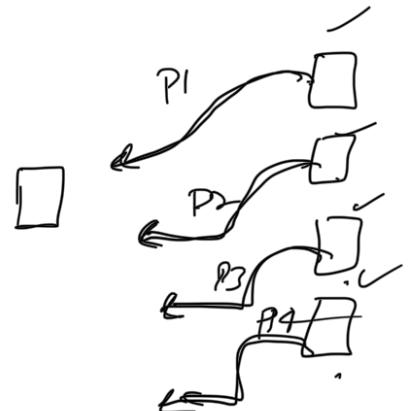
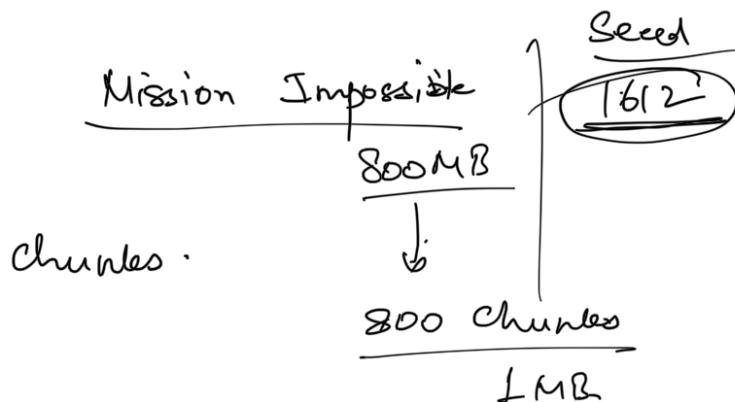
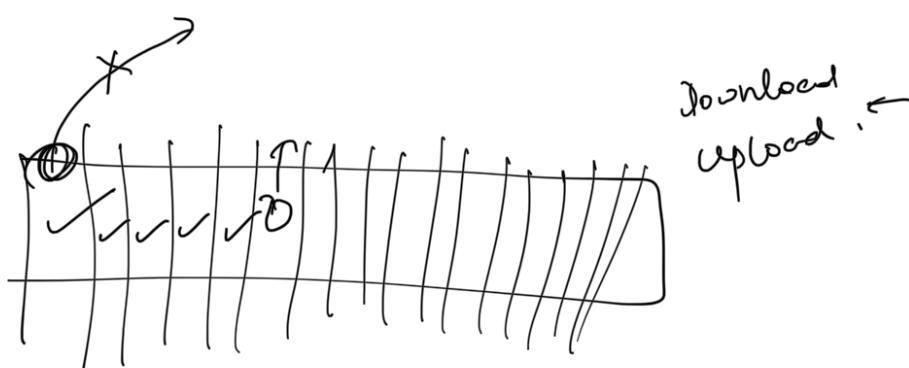


① Storage should be able to store large files

② Reliable & durable

③ Retry / Network failure → gracefully work repeatable
Should be minimized

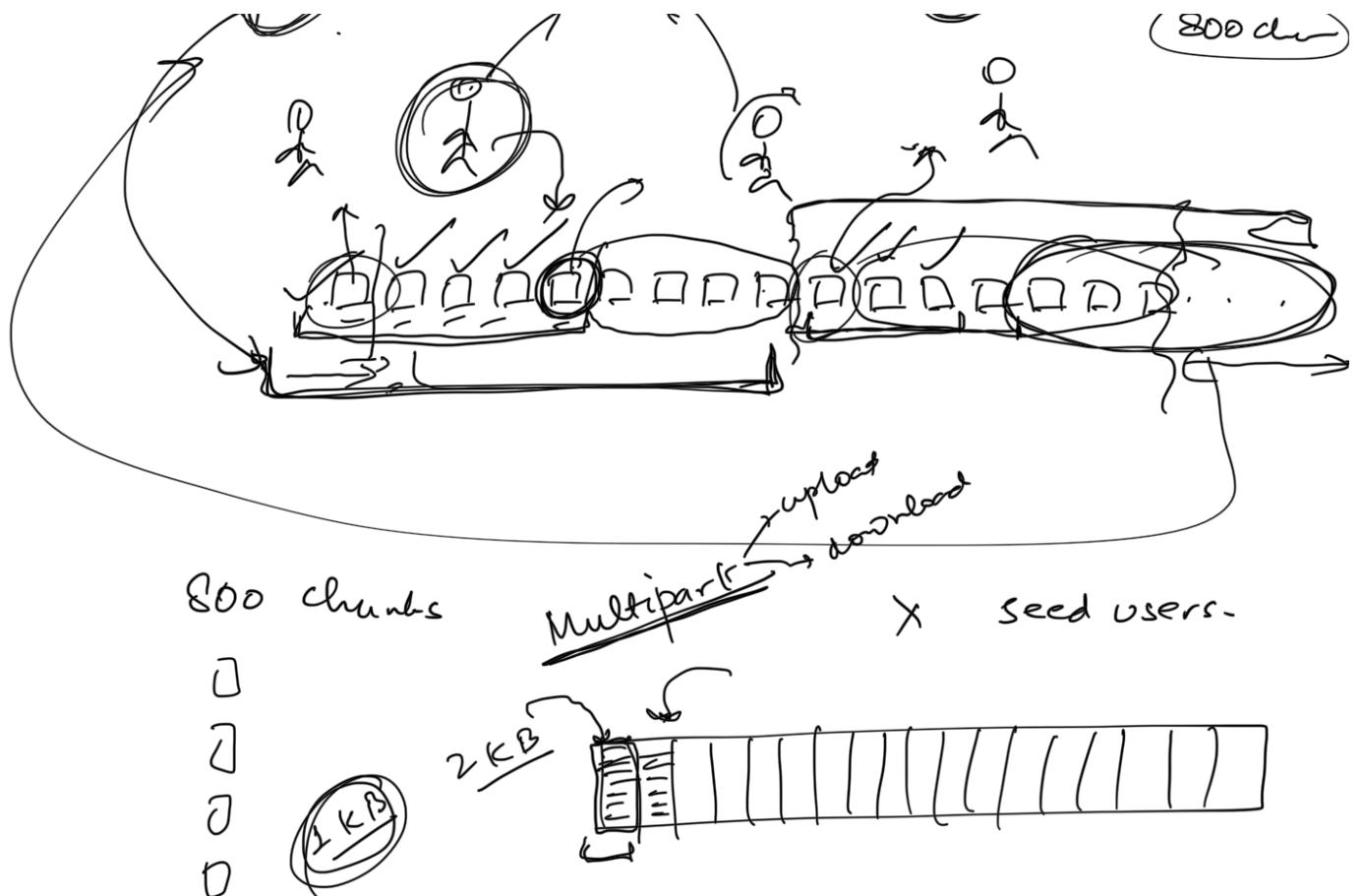
④ Download the uploaded file



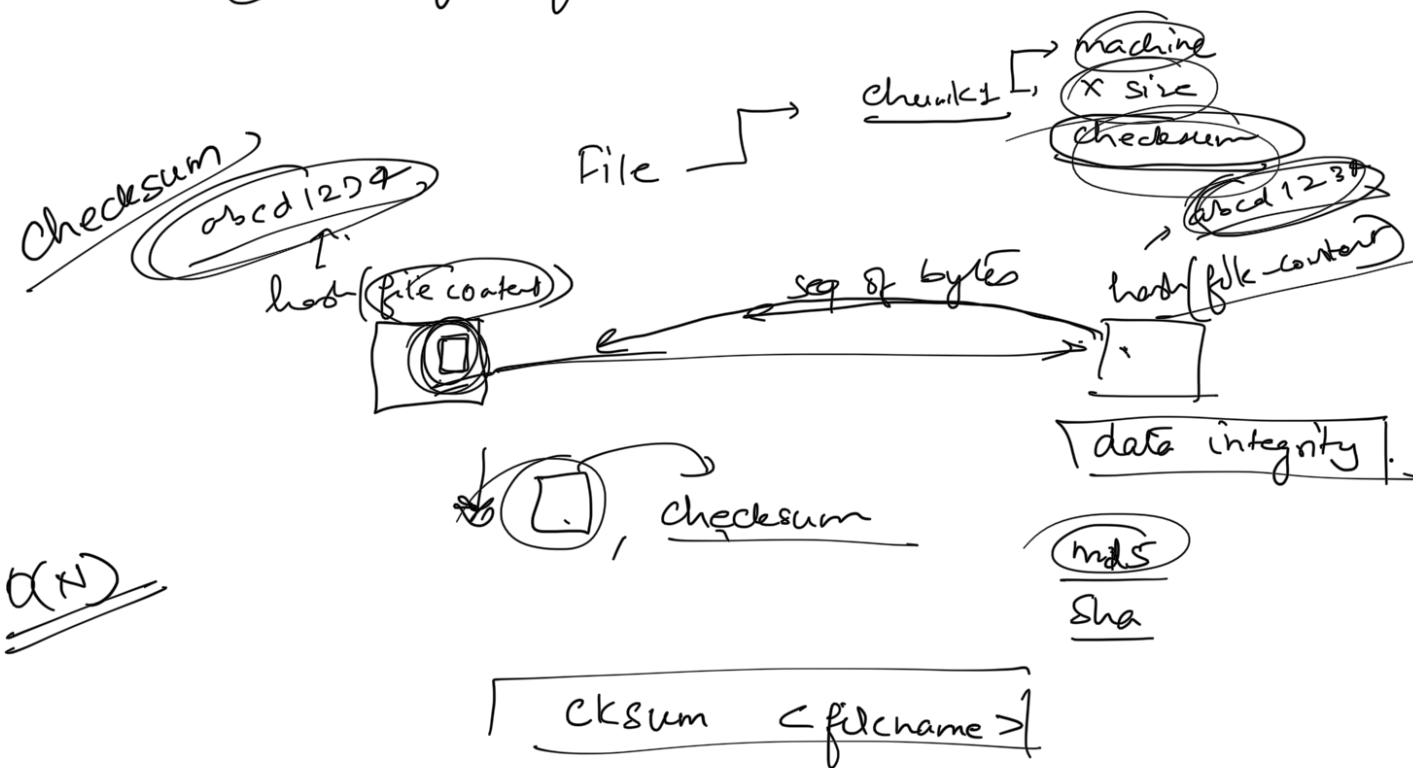
mil·ani

seed users





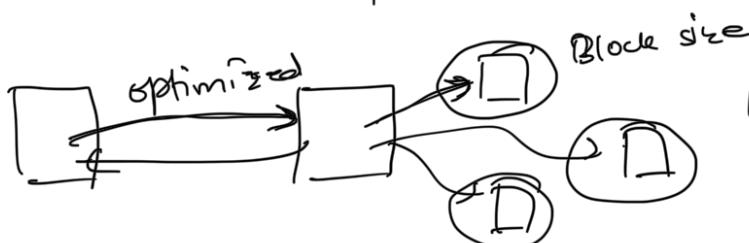
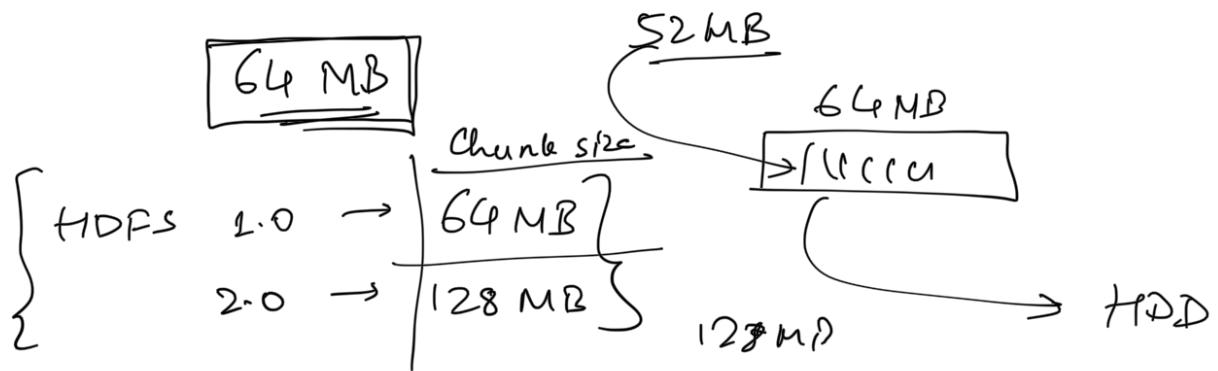
① Large file → break into chunks.



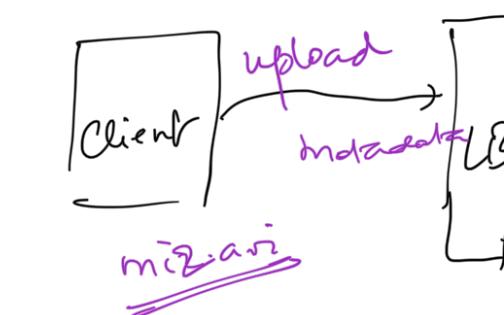
1 GB

Total size → 4 MB → default HTTP req limit 500 MB

① 4 MB → default HTTP req limit



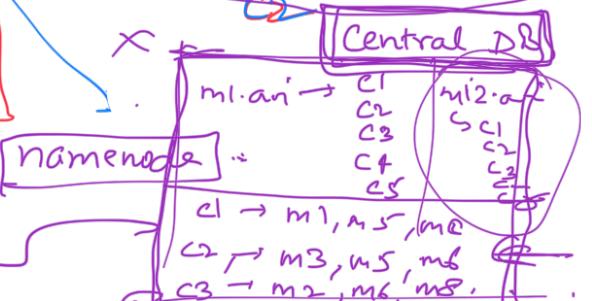
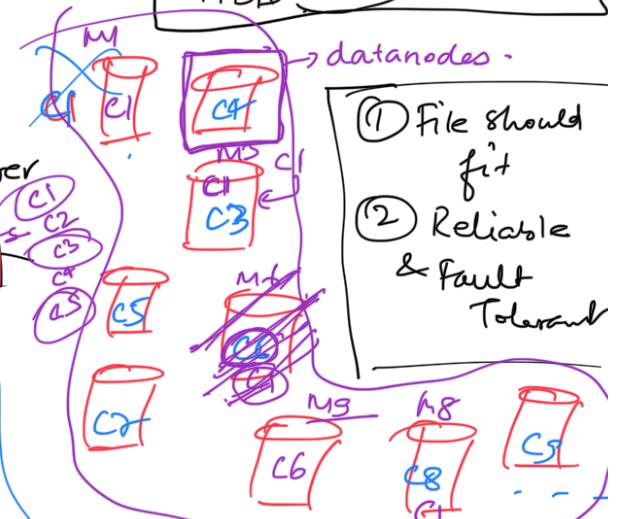
C1 66



2 copies
2 copies

500 GB → 1000 chunks

File → chunk
all chunks → locations



F, chunk → (m1, m2, m3)

m1

m3

m5

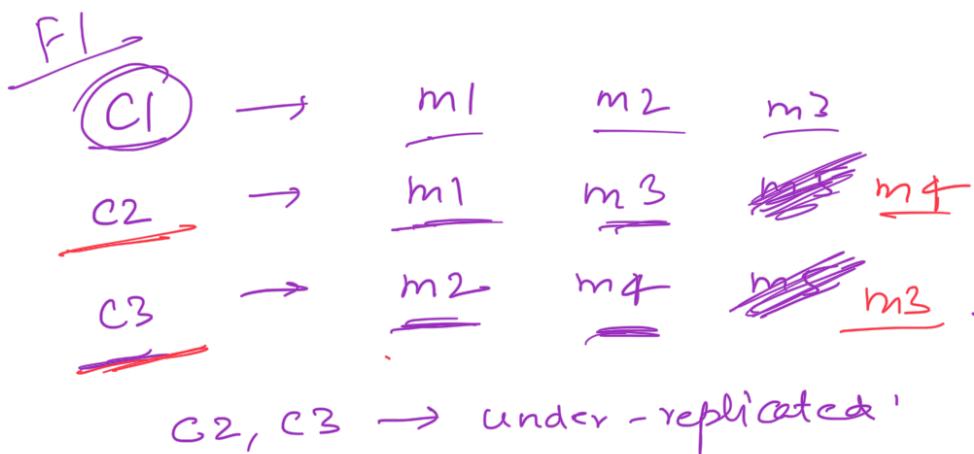
Under-replicated

...

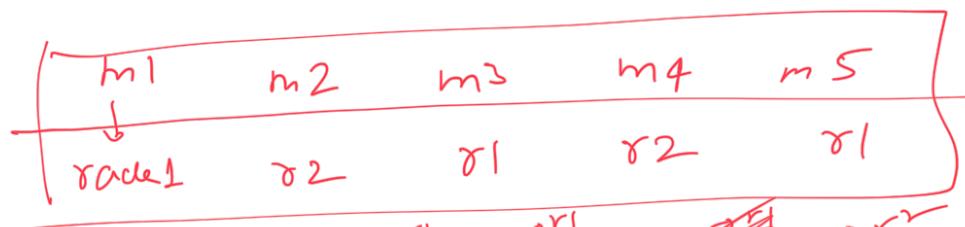
Rack-aware



2



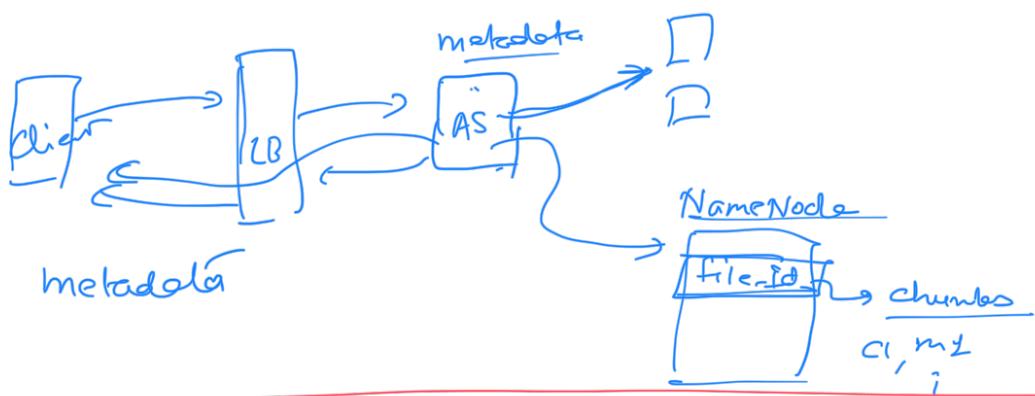
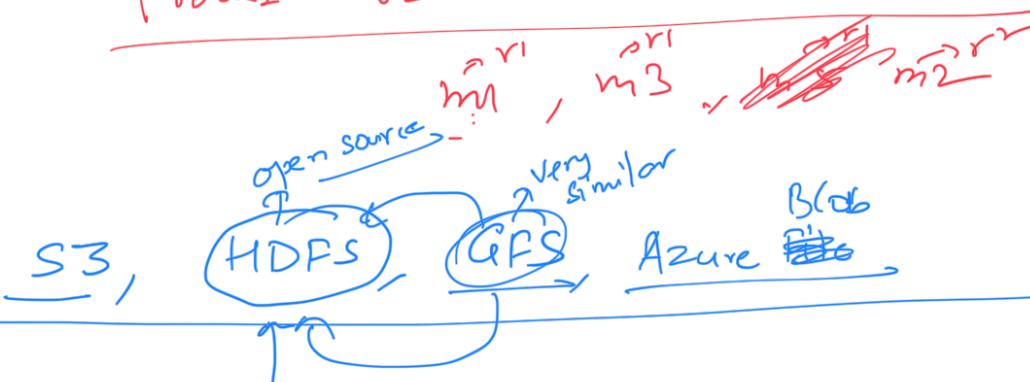
replication-level ≈ 2



if (rack-aware){

check if at
3 on 5

if (yes){
do}



PROBLEM #2:

Nearest neighbor problem

(lat, long)

→ find nearest ~~to~~ restaurants

Brute Force get-closest-X(lat, long):
 for rest in all-rests:
 $\text{dist} = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$
 if dist of markheap.pop() is
 greater than dist:
markheap.pop
markheap.insert(vot)

Very very slow

SQL

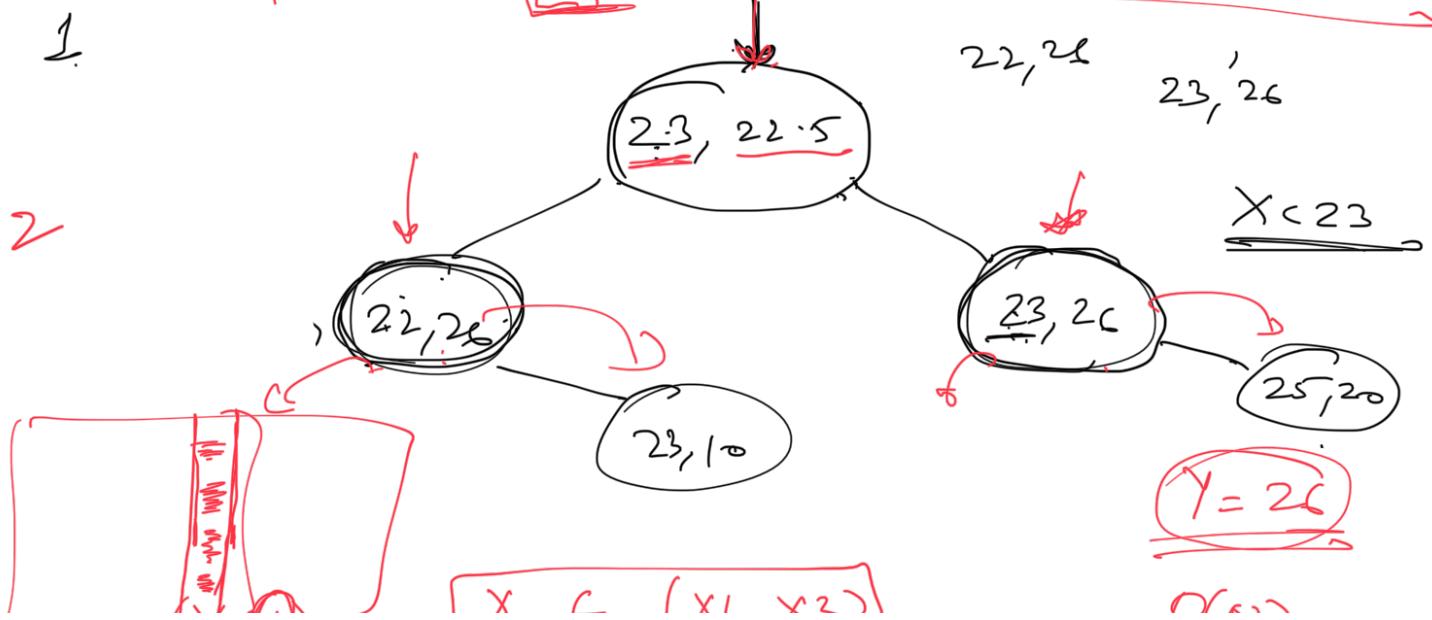


rest_id	name	X	Y
		X	Y

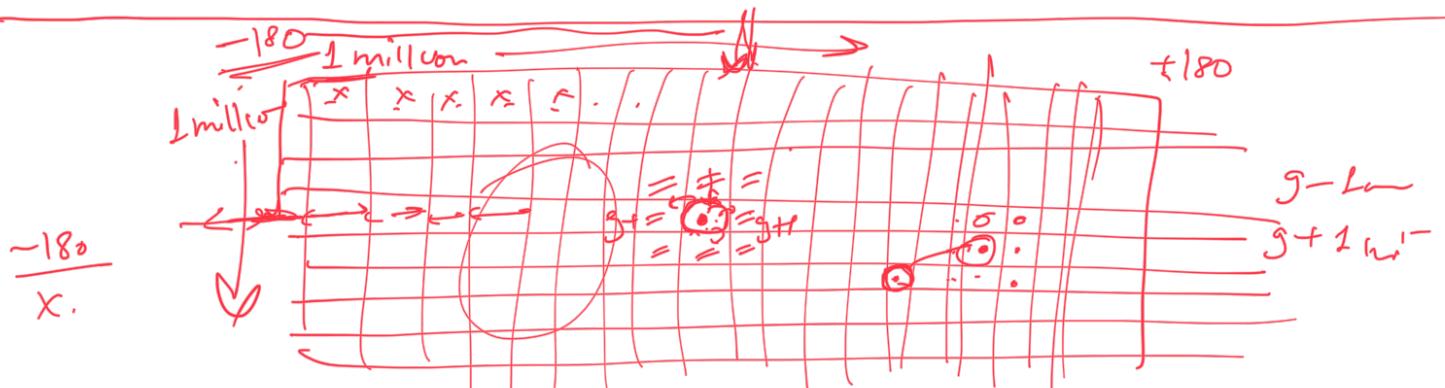
index(X, 4)

• (lat, long) k know SELECT * FROM restaurants
 WHERE
 $\rightarrow X \text{ between } (\text{lat}+k, \text{lat}-k)$
 $\rightarrow Y \text{ between } (\text{long}+k, \text{long}-k)$

Approach 2



~~(x_1, y_1)~~ → $\exists x \in (x_1, x_2)$ ~~y_2~~ $\exists y \in (y_1, y_2)$ ex



(lat, long) → grid-id
 $-180 + x * g = \text{lat}$

rest → grid

$$y = \frac{\text{lat} - 180}{g}$$

indexed

rest_id	name	lat	long	grid-id

grid-id
 \downarrow

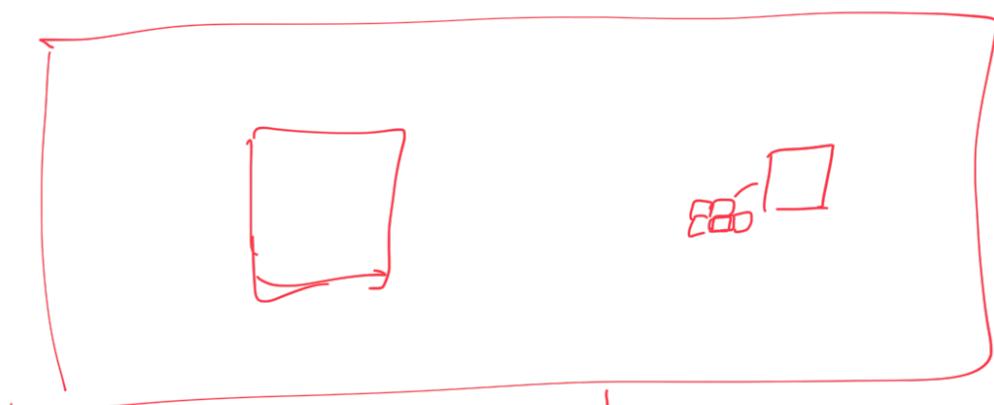
neighboring grid-ids

FAST

enough results.

SELECT * FROM restaurants

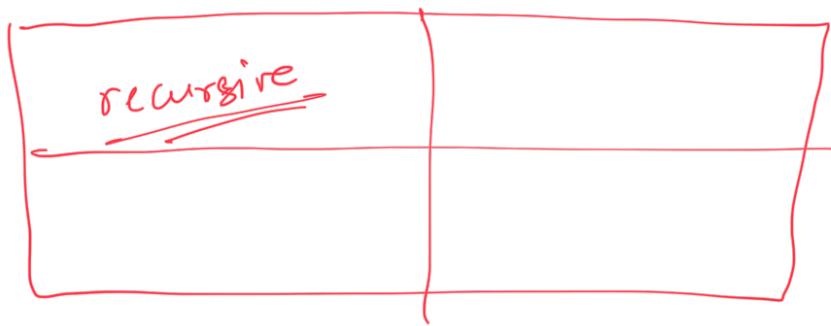
WHERE grid-id IN (g_1, g_2, g_3, \dots)



~ 50 rest

- ① Variable size grids
- ② $(x, y) \rightarrow \text{grid}$
 \downarrow
 neighboring grids.

①



□ CSD
rest

