

Agenda

pragy@scaler.com

- Working with large files
- Nearest Neighbour Queries
Quad Tree



Images / Videos → user generated content ⇒ few Mb's }
few Gbs. }

App server Logs → append only files → 50 Tbs
100 Tbs

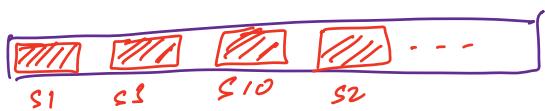
Stock broker / Financial Institute

↳ mandated by govt to maintain
log-

Requirements

- store extremely large files (100 TB)
- storage must be reliable & durable
- Both read & write must be possible
 - ↳ read → download file on demand stream
 - perform analytics on top of the file.

How to store large file → split it into **chunks**
& distribute chunks across servers



Chunk Size → 128 Mb in HDFS

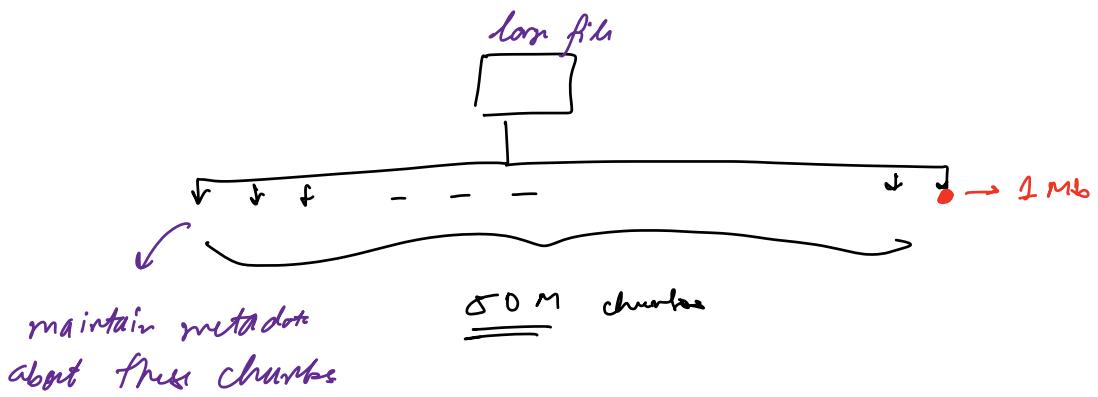


1 Mb file of 50 TB

of chunks

$$\frac{50 \times 10^{12} \text{ bytes}}{1 \times 10^6 \text{ bytes}} = 5 \times 10^7$$

= 50 Million chunks



- ① Merging/collation → if client wants to download we have to merge these chunks.

↓
overhead grows as chunk size gets smaller.



System 32 objects in Windows
node-modules/ - - -

\rightarrow Same 3Gb in total.
 1 Mb/sec

- ② Store metadata about chunks

every chunk → ① fill it is part of → 8 bytes

② offset → 8 bytes

③ server stored on. → 4 bytes

20 bytes / chunk.

$$50 \text{ M} \times 20 \text{ bytes} = 1000 \text{ M bytes} = 1 \text{ GB}$$

of metadata.

HDFS

① Name Node

→ also replicated (Master-slave)

immediate consistency

↳ only 1 name node

↳ stores file metadata

② Data Node

↳ stores chunks

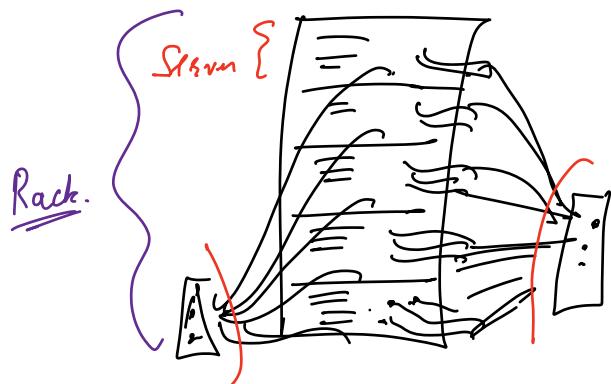
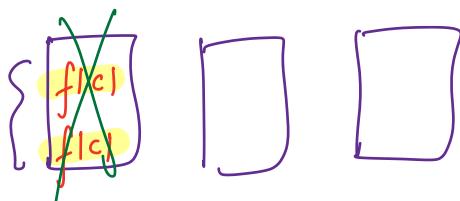
every chunk is replicated

& stored in multiple data nodes



↳ rack-aware algorithms

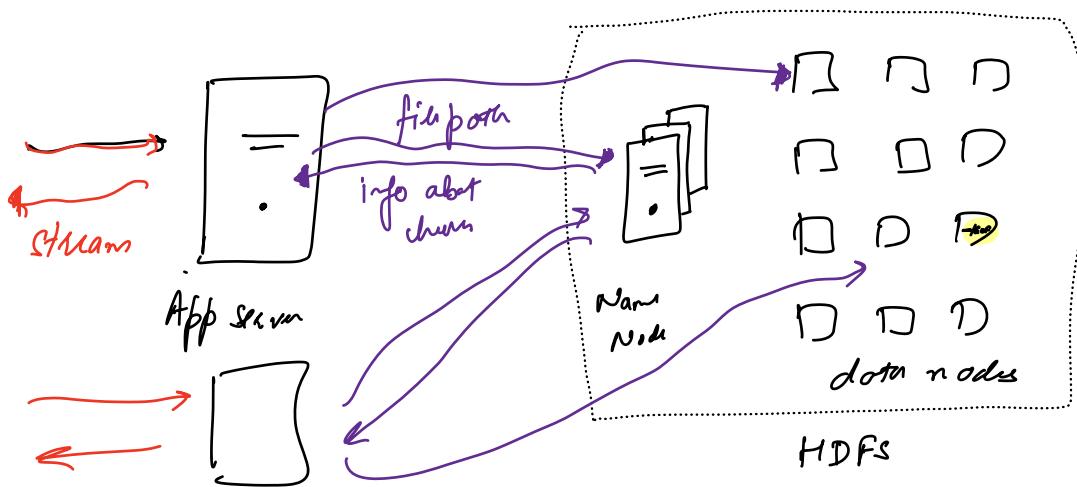
replicas
should be
on same
machine



Name Node

filepath	chunk id	Server-side
log-all.txt	0	[10, 5, 3]
log-all.txt	1	[7, 2, 1]
log-all.txt	2	[1, 3, 7]
:		

Read / Write in HDFS



① Read

- ↳ go to nam node & find metadata
 - ② dont file's chunks
- ↳ for each chunk (sequentially in order of chunks)
 - ↳ gets servers containing chunk
 - ↳ load the chunk in mem (128 MB)
 - ↳ Stream to the client

Map-Reduce
aggregate function

② Write

↳ client will start sending large data stream to app serv.

↳ app serv will go to NameNode & find chunk size

↳ for now app serv get 128MB in buffer

① go to NameNode

where should I store
this new

chunk?

← seen id

② go to server-id &
contt this 128MB
chunk

③ NameNode → I have
written it this chunk-id
at server-id

→ clear in buffer & continue receive
from client

App Server

```
byte [128*1000+1000] buffer; → in-memory  
int count = 0;
```

```
void receive () {  
    byte data;  
    while (data = receivefromClient() ≠ EOF) {
```

```
        buffer [count++] = data;
```

```
        if (count == 128*1000+1000) {
```

```
            store (buffer) //async
```

①

}

```
        count = 0;
```

}

{

What if a file is not in multiples of 128MB

5MB

{

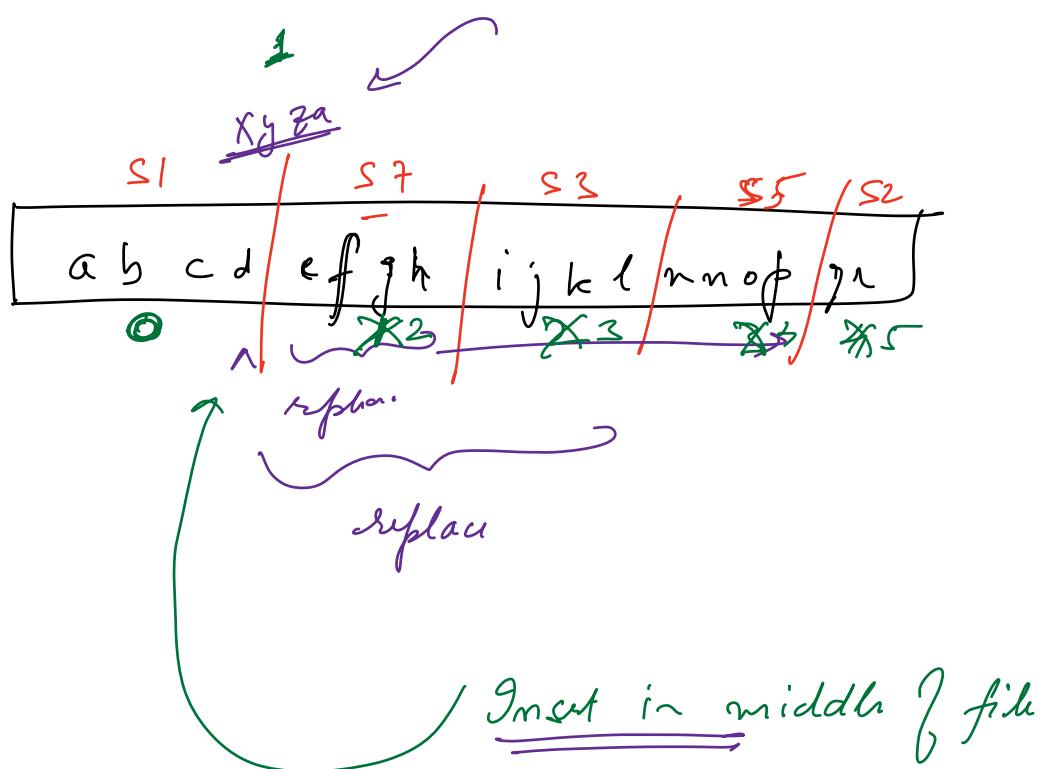
200MB

{

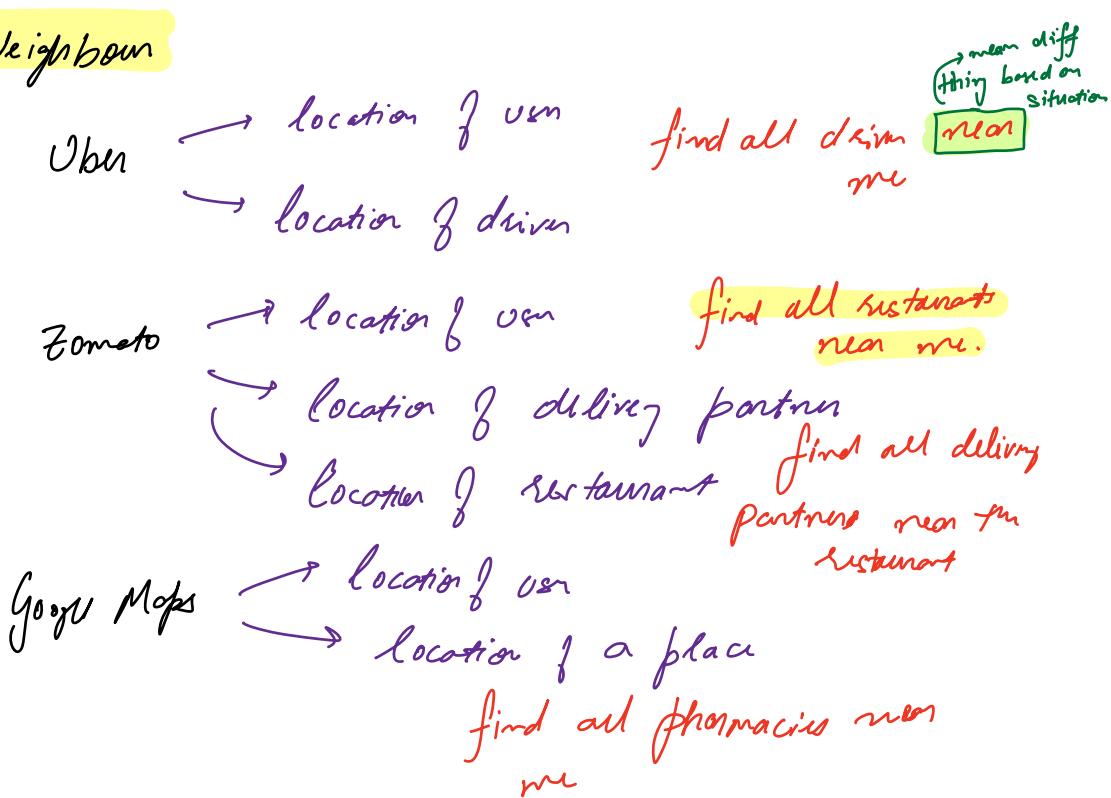
↓
 1 cherr
 $C_0 \rightarrow 5MB$
 ↓
 2 cherr
 $C_0 \rightarrow 12MB$
 $C_1 \rightarrow 72MB$

file part

10:21 → 10:35



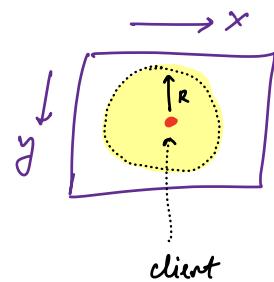
Nearst Neighbour



SOL

Restaurants

id name -- lat lon



Circle Query

query for find all restaurants near me.

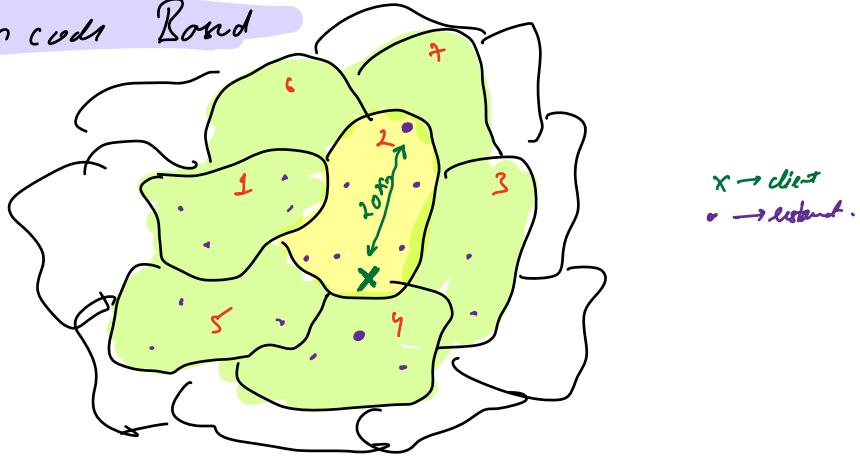
Select * from restaurants
where euclid-distance-2d(lat, lon, user-lat, user-lon) $\leq R$;

Suppose radius = R
contin fab
linear scan.

euclid-distance-2d(x_1, y_1, x_2, y_2) = $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

1.4 million restaurants

Pin code Based



Pro → easy to query.

resturants

id | name --> pincode.

select * from rest

where pincode = User-pincode

Con → other resturants in
nearby pincodes shot
are closer.

Pincode neighbors

Pincode	neib ~ P^i
2	1
2	3

→ what if a pin code
is large / neighbors
distance

→ what if there are too
many resturants in the
pincode?

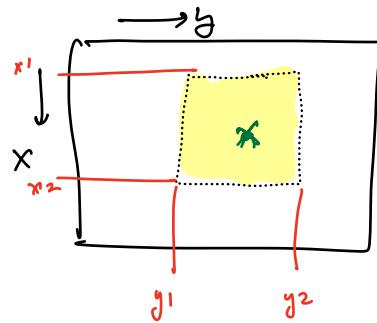
slow queries

Rectangular Box Based

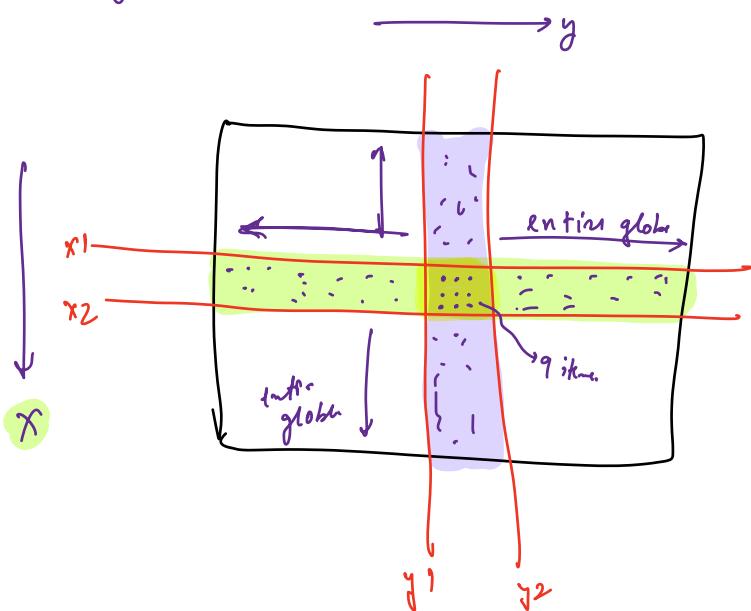
slow
but fast
not circle

Select * from restaurant
where lat between x_1 and x_2
and lon between y_1 and y_2

index index



index scan on lat
linear scan on lon



A) find via index

linear search for $y_1 \leq \text{lon} \leq y_2$

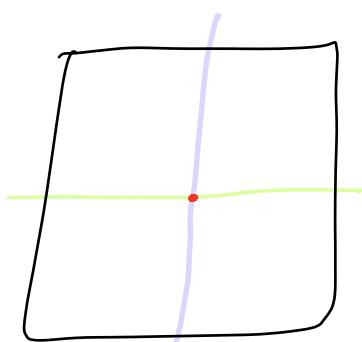
B) find via index

linear scan for $x_1 \leq x \leq x_2$

C) find via index

find via index

intersect.



K-d tree

→ generalization of DS that allows you to use multiple indices at the same time.

bad worst case $O(N)$

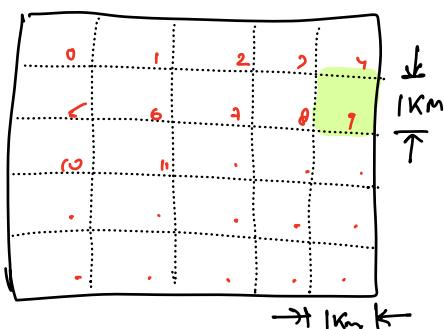
v. good avg case $O(\lg N)$

→ it builds a k -dimensional index

K-d tree $K=1 \rightarrow \text{RST}$
K-d tree $K=2 \rightarrow \text{Quad Tree}$
K-d tree $K=? \rightarrow \text{Oct Tree}$

$k \ll N$ (if $k \approx n \rightarrow$ curse of dimensionality)
 dimensionality. # of items (drain location, rest location)

Grid Based → basically pin-code on steroids
 we control for shape
 we control for size



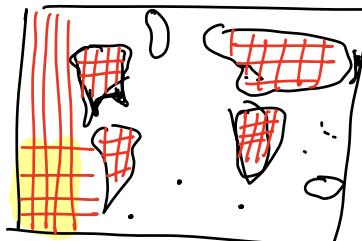
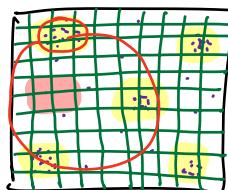
rect
id / raw / cell-id

Challenges

- ↳ high density hubs
- ↳ low density areas

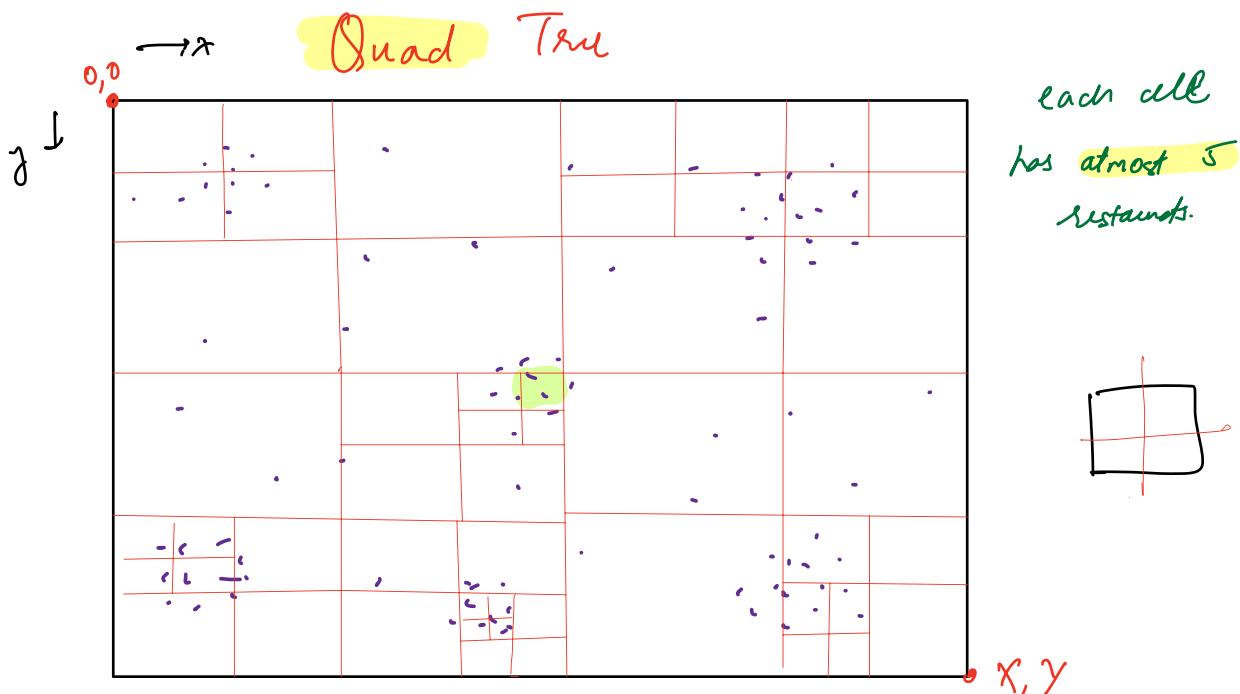
No optimal grid size

- ↳ high density
 - ↳ search smaller neighborhood
- ↳ low density
 - ↳ expand search to wider area.



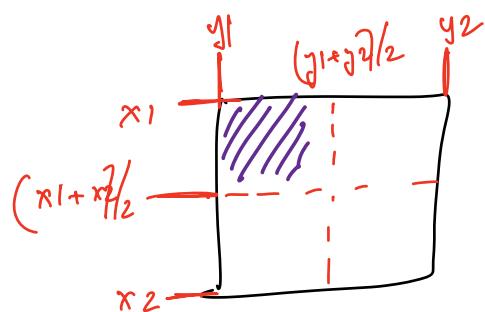
Magic Solution

→ a dynamic grid that auto
adjusts to density of a region.



start with the entire world as a node

- ↳ if node has > 5 restaurants
 - ↳ split into 4 quadrants of equal size
 - ↳ recurs on each quadrant.



dimensions of current cell

```

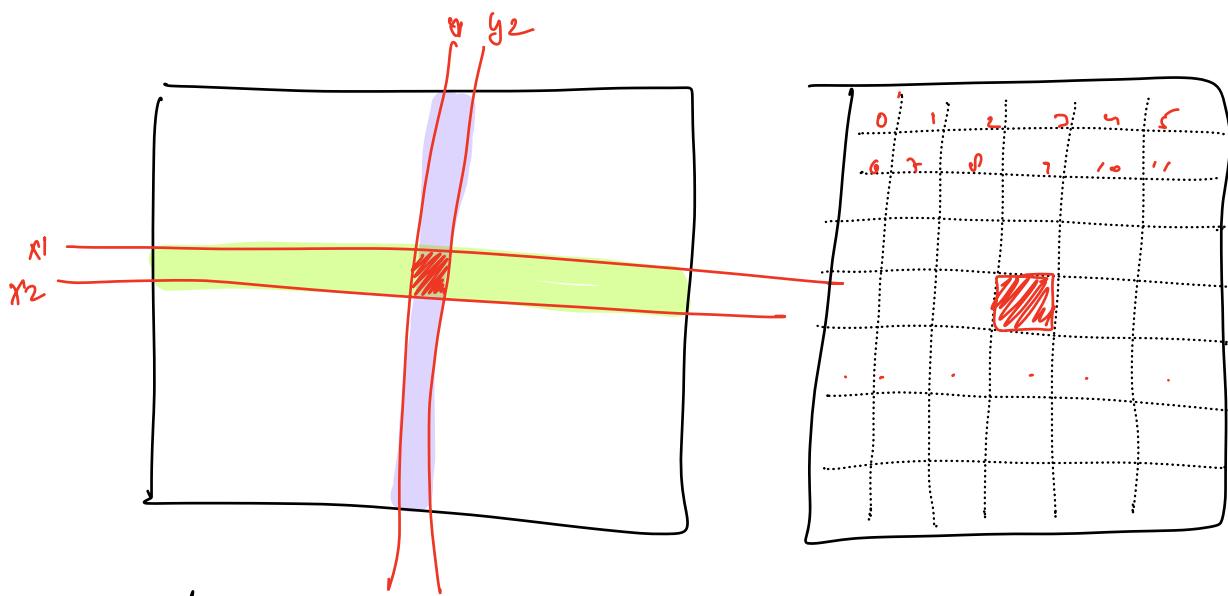
TruNode create-grid ( List< Point2D> rectangles, x1, x2, y1, y2 ) {
    if (rect. extent <= 5)
        return
    TruNode tl = new TruNode ( ≡, x1, xmid,
                                y1, ymid )
    TruNode tr = new TruNode ( ≡, x1, xmid )
                                ymid, y2 )
    TruNode bl = new TruNode ( ≡, xmid, x2 )
                                y1, ymid )
    TruNode br = new TruNode ( ≡, xmid, x2 )
                                ymid, y2 )
}

```

upload to GDrive \rightarrow 10Mbps

\hookrightarrow 200kbps for small file

Zip file \rightarrow 300Mbps
10Mbps



id/nan/x/y

id/nan/cell-id

sel + —

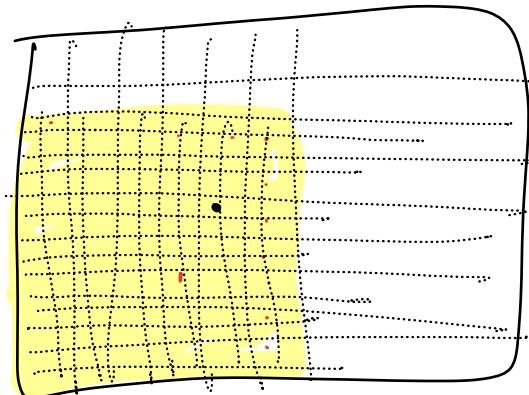
wh x let $x_1 \Delta x_2$
y in $y_1 \Delta y_2$

wh + —

wh all id = —



18 alle ✓



68 alle

pragy@scaler.com.