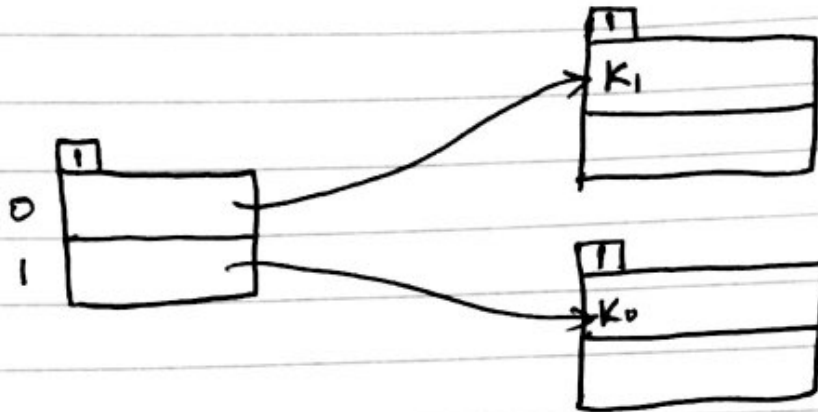
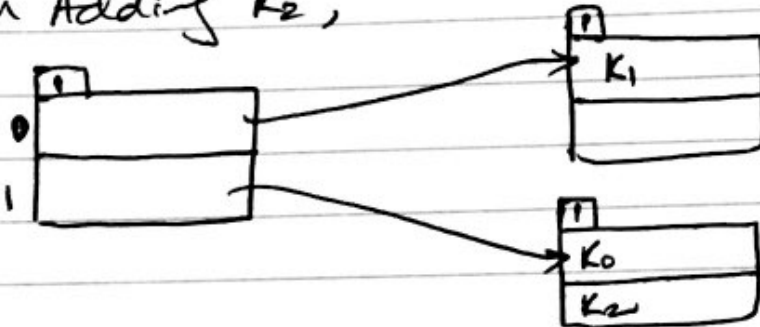


HOMWORK 4

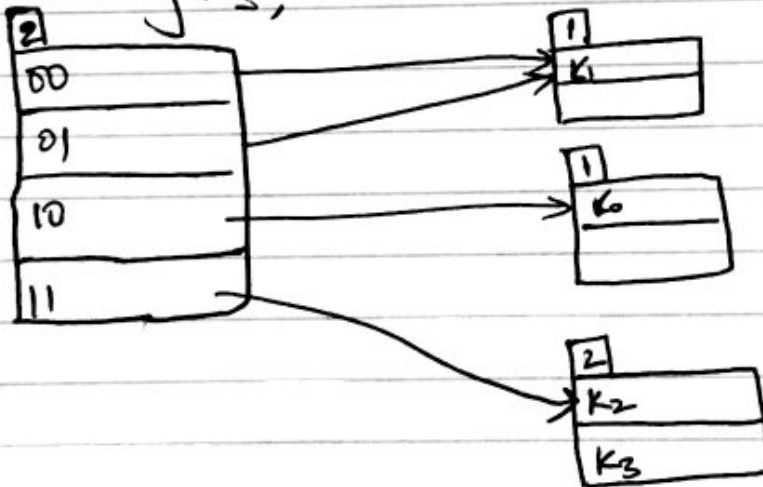
Q.1.)



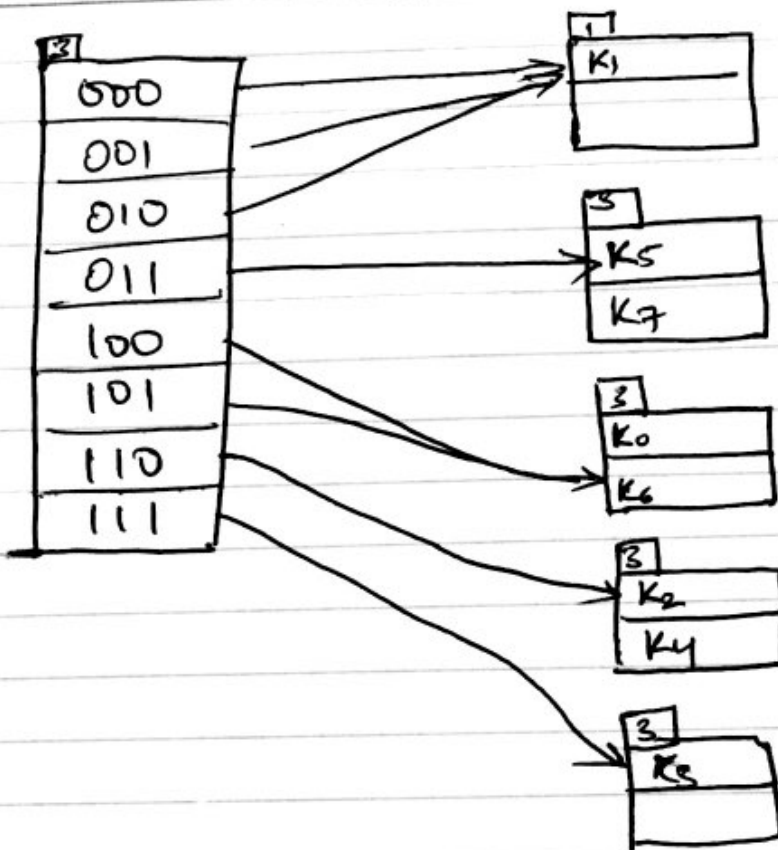
On Adding K₂,



On Adding K₃,



Similarly on adding K₄, we have to increase the keys further in the table.



Q₂)

$$\text{RPM} = 150,000$$

$$\text{Tracks} = 250,000$$

$$2000 - \text{Sectors}$$

$$\text{Sector Size} = 512 \text{ bytes}$$

$$\begin{aligned} \text{Capacity of disk} &= 2 \times 2 \times 250,000 \times 2000 \times 512 \\ &= 1.024 \text{ Tb} \end{aligned}$$

$$\text{Size of each track} = 2000 \times 512 = 1.024 \text{ mb}$$

So, 1mb can be read in 1 rotation

$$\begin{aligned} \therefore \text{Maximum rate at which data can be read from disk} &= \frac{15000 \times 1.024}{60} = 256 \text{ mb/s} \end{aligned}$$

$$\text{Average rotational latency} = 0.5 \times 1000 \times \frac{60}{15000}$$

$$= 2 \text{ ms}$$

(b) Seek Time = 5 ms
Block Size = 4 kb

$$\text{Transfer Rate} = \frac{4 \text{ kb}}{250 \text{ mb/s}} = 0.016 \text{ ms}$$

$$\text{Time to Read} = \text{Seek Time} + \text{Rotational latency} + \text{Transfer Rate}$$

$$= 5 + 2 + 0.016$$

$$= 7.016 \text{ ms}$$

BLOCK MODEL

For file size of 100 Kb ,

$$\text{Blocks} = \frac{100 \text{ Kb}}{4 \text{ Kb}} = 25$$

$$\therefore \text{Transfer Rate} = 7.016 \times 25 = \underline{175.4 \text{ ms}}$$

For file size of 1000 Kb ,

$$\text{Blocks} = \frac{1000 \text{ Kb}}{4 \text{ Kb}} = 250$$

$$\therefore \text{Transfer Rate} = 7.016 \times 250 = \underline{1754 \text{ ms}}$$

for file size of 100 Mb,

$$\text{Blocks} = \frac{100 \text{ Mb}}{4 \text{ Kb}} = 25,000$$

$$\therefore \text{Transfer Rate} = 7.016 \times 25,000 = 175,400 \text{ ms}$$

Latency / Transfer Model

for file of size 100 Kb,

$$t_r = 2 + 5 + \frac{100}{256} = 7.3 \text{ ms}$$

for file of size 1000 Kb,

$$t_r = 2 + 5 + \frac{1000}{256} = 10.9 \text{ ms}$$

for file of size 100 mb,

$$t_r = 2 + 5 + \frac{100,000}{256} = 397.625 \text{ ms}$$

(c)

File Size = 36 GB

Main Memory = 1 GB

We can do 36 files of size 1 GB each
= $36 \times (\text{Time to sort 1 file})$

$$= 36 \times \left[5 + 2 + \frac{1 \times 10^3}{256} \right] \times 2 = \underline{\underline{281.85}}$$

Considering 36 splits ;

$$\text{Buffers} = 36 + 1 \text{ (o/p buffer)} \\ = 37$$

$$\text{Buffer Size} = \frac{1 \text{ GB}}{37} = 27.02 \text{ mb}$$

$$\text{Total files} = \frac{36}{27.02 \text{ mb}} = 1333 \text{ files}$$

$$\text{Time to read \& write these files} = 1333 \times \left[\frac{5+2+\frac{27.02}{256}}{2} \right] \\ = 306.8 \text{ s}$$

$$\text{Total Time} = 306.8 + 281.75 \text{ s} \\ = 588.55 \text{ s}$$

(d) $d=6 \Rightarrow 2 \text{ phases.}$

$$\text{Buffer} = 6(\text{input}) + 1(\text{o/p}) = 7 \\ \text{Buffer Size} = \frac{1 \text{ GB}}{7} = 142.86 \text{ MB}$$

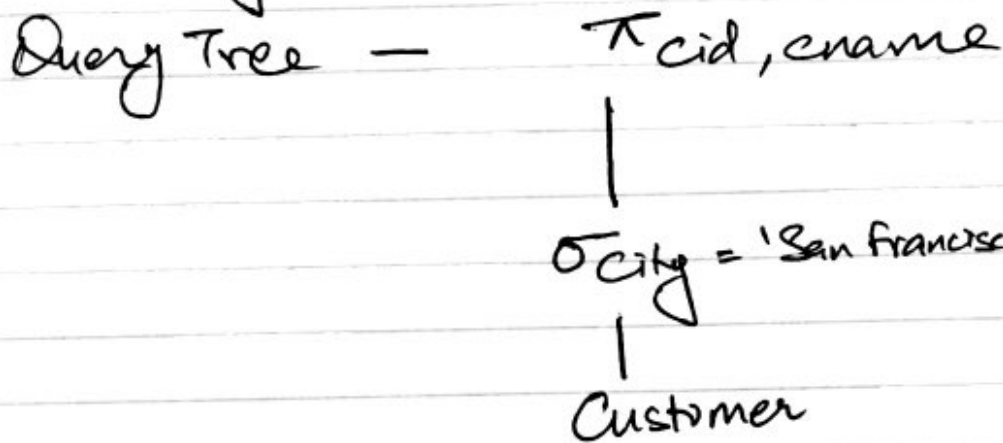
$$\text{Total files} = \frac{36 \text{ GB}}{142.86 \text{ MB}} = 252$$

$$\text{Time to read \& write these files} = 252 \left[\frac{5+2+\frac{142.86}{256}}{2} \right] \\ = 284.78 \text{ s}$$

$$\text{Total Time} = 284.78 \times 2 = \underline{569.56 \text{ s}}$$

Q3)

(a) Query 1



Total number of customers in San Francisco = 1 million

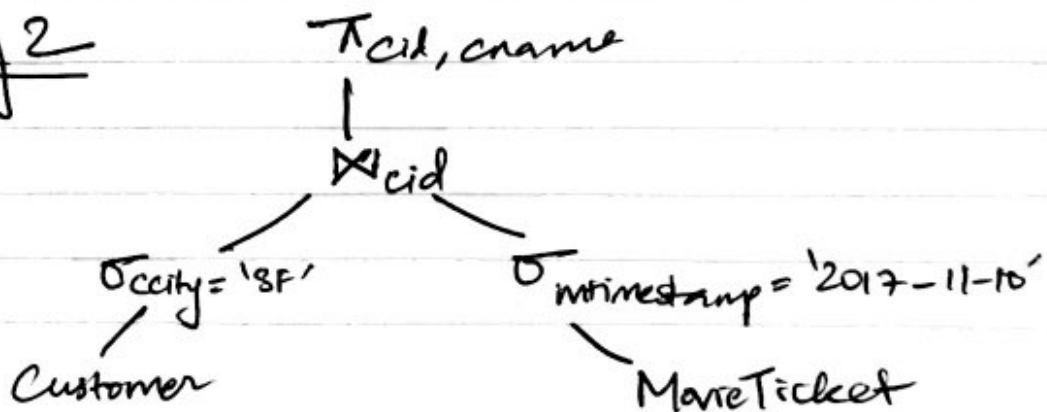
Main Memory size = 1.5 GB

Main memory data = 1 million \times 100 bytes (tuple size) = 100 MB

Now, because there is no index we don't know when to stop looking for customers = $100 \text{ MB} \times 100 \text{ bytes} = 10 \text{ GB}$

$$\text{Time to read} = 10 \text{ ms} + \frac{10 \text{ GB}}{150 \text{ MB/s}} = 66.67 \text{ s}$$

Query 2



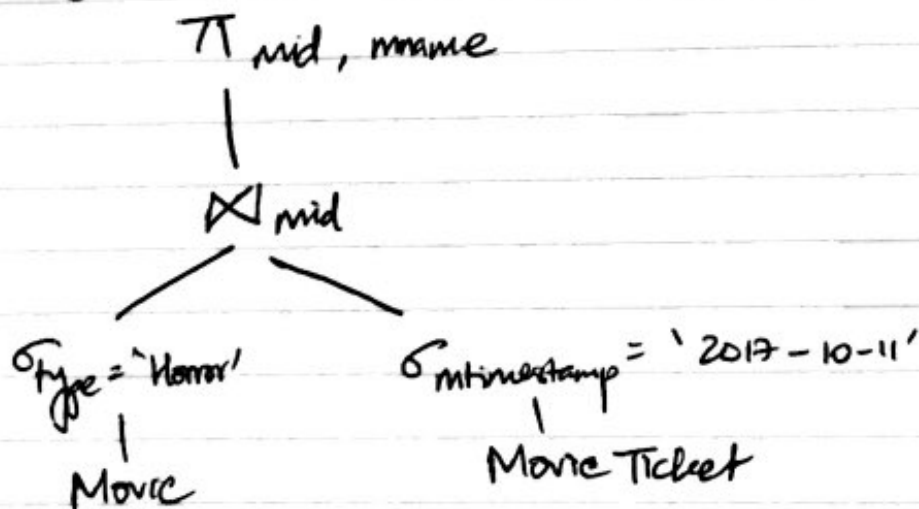
Here we join 2 tables customer and movieticket where city is San Francisco and date 2017-11-10.

To scan for SF take 66.67s
Movie Ticket has 5 billion rows of 40 bytes.

$$\text{Cost} = 10\text{ms} + \frac{5 \text{ billion} \times 40 \text{ bytes}}{180 \text{ mb/s}} = 1333.33 \text{ s}$$

$$\text{Total Time} = 66.67 + 1333.33 = \underline{\underline{1400 \text{ s}}}$$

Query 3



$$\begin{aligned} \text{Cost} &= \left(10\text{ms} + \frac{0.5 \text{ million} \times 100 \text{ bytes}}{180 \text{ mb/s}} \right) + \left(10\text{ms} + \frac{5 \text{ billion} \times 40 \text{ bytes}}{180 \text{ mb/s}} \right) \\ &= 0.34 \text{ s} + 1333.33 \text{ s} \\ &= \underline{\underline{1333.67 \text{ s}}} \end{aligned}$$

(6) City size = 32 bytes
CID size = 8 bytes
KID = 8 bytes

Each Tree node has $n-1$ keys & n -pointers.

Since Customer table has a sparse index,

$$(n-1) * 32 + n * 8 = 4096$$

$$n = 103$$

Assuming occupancy of 80%, we will have 82 indexes.

Considering Customer table = 100 million

Entry size Customer = 100 bytes.

Records per disk page = $4096 / 100 = 40$ records

No of indexes = $\frac{100 \text{ million}}{40} = 2.5 \text{ million}$

We will have 2.5 million blocks of 40 records.

$$\text{No of leaf nodes} = \frac{2,500,000}{82} = 30488 \text{ nodes}$$

$$\text{No of nodes in next level} = \frac{30488}{82} = 372 \text{ nodes}$$

$$\text{No of nodes in next level} = \frac{372}{82} = 4 \text{ nodes}$$

So, the B+tree will have 4 levels of nodes: the root, two internal levels and leaf level.

$$\text{Size of tree} = 30480 \times 4\text{Kb} = 121 \text{ Mb}$$

$$\text{Time to fetch a single record} = 5 \times 10 \text{ ms} = 50 \text{ ms}$$

$$\text{Time to fetch 50 records} = (4+3) \times 10 = 70 \text{ ms}$$

DENSE UNCLUSTERED INDEX ON (CID):

$$\begin{aligned} \text{Each Index Entry} &= 8(\text{CID}) + 8(\text{RID}) \\ &= 16 \text{ bytes} \end{aligned}$$

$$\text{Nodes} = \frac{4096}{16} = 256$$

Assuming 80% occupancy, we will have 204 entries per node.

Movie-Tickets has 5 billion records.

$$\text{No of leaf nodes} = \frac{5 \text{ billion}}{204} = 24.5 \text{ million}$$

$$\text{No of nodes in next level} = \frac{24.5}{204} = 120000$$

$$\text{No of nodes in next level} = \frac{120000}{204} = 588$$

$$\text{No of nodes in next level} = \frac{588}{204} = 2 \text{ nodes}$$

~~next level~~

So, the B+ tree will have 5 levels.

$$\text{Size of tree} = 24.5 \text{ million} \times 4 \text{ KB} = 98 \text{ GB}$$

$$\text{Time to fetch a single record} = 6 \times 10 \text{ ms} \\ = 60 \text{ ms}$$

$$\text{Time to fetch 50 records} = (49 \times 10 \text{ ms}) \\ + 60 \text{ ms} \\ = 550 \text{ ms}$$

(C)

Query 1 - I would choose a clustered index on cid and ccity in customers, to accelerate the process of fetching records from customers.
To fetch 1% of customers from SF it would take 0.1 sec.

Query 2 - Choosing clustered indexes on ccity in customers and in timestamp in Movie Ticket would be ideal. We would first fetch users in SF which would take only 1% of the cost and then further filter for the date. Therefore the cost would be reduced further.

Query 3 - We choose unclustered indexes on ~~film~~ in ~~Reviews~~ ^{Movie} and clustered on Movie Ticket. The first index can fetch users from SF, which would take only 1% of the cost. Then we can scan Movie Ticket for timestamp '2017-10-11', hence improving runtime.