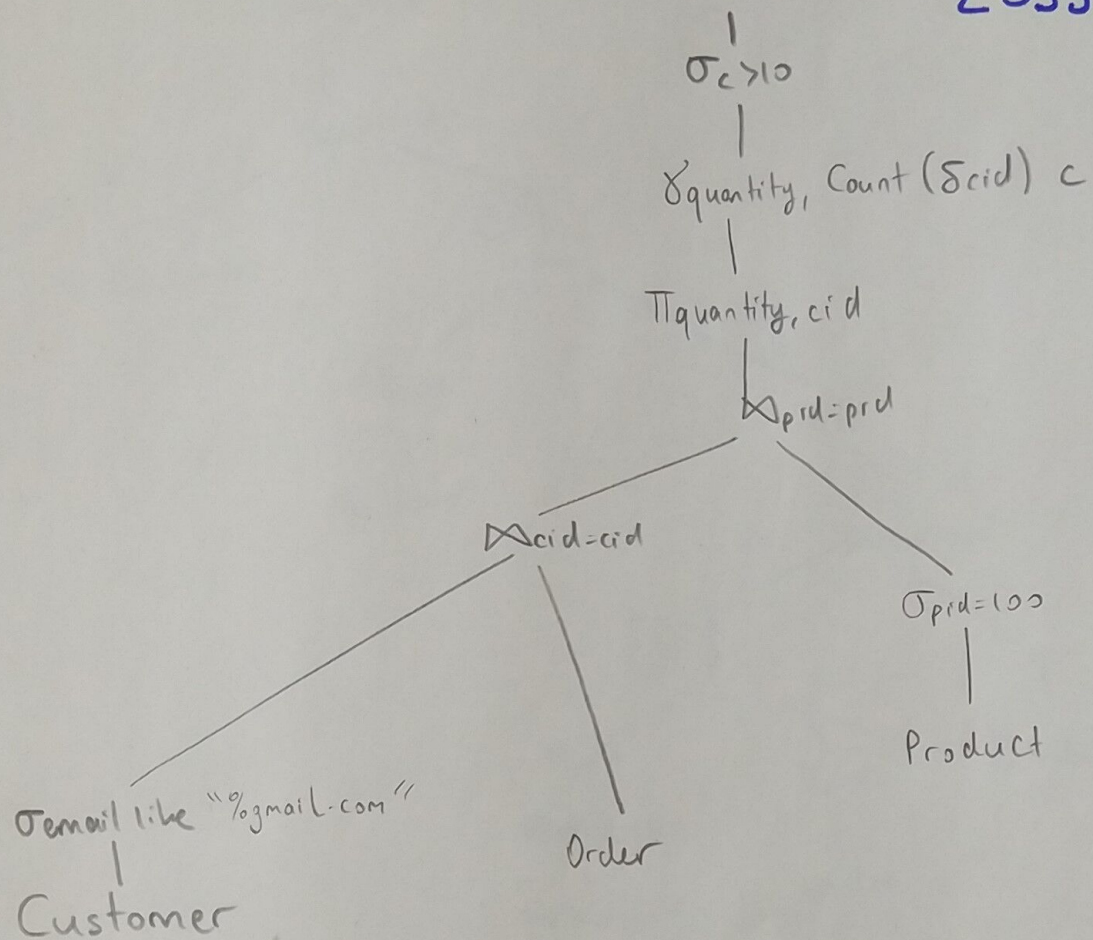
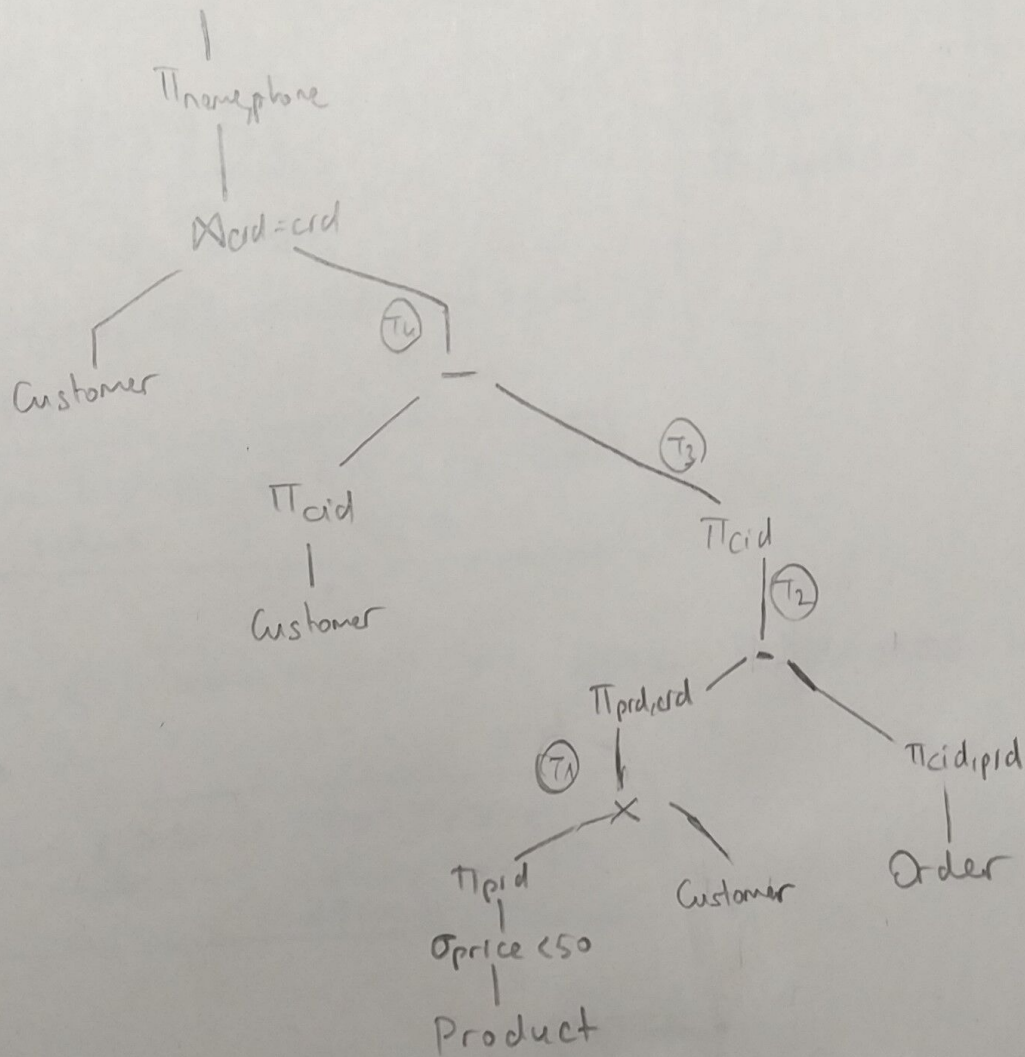


1A)





1B)



Clarification:

$T_1$  = We need all customer and product cross product for all possible prod-cid pairs

In the left, we have all pairs with price less than 50, in the right we have all orders. The difference will give us orders that actually ~~does not~~ exist (with less than 50 price)

Then, project cid in  $T_2 \rightarrow T_3$

outer query

Select cid  
From Customer  
Except

Select cid  
From  $T_3$  (or  $T_3$ , they are the same)

=  $T_4$

In the final step

(Select name, phone  
From Customer,  $T_4$   
Where  $T_4.cid = Customer.cid$ )



## 2A)

i) Table scan

$$B(\text{Teach}) = \boxed{500 \text{ I/O.}}$$

ii) Clustered B+ on semester

- 4 different values. each semester will have approx 125 page.

$$\text{Cost} = \boxed{\text{Some I/O for index} + 125 \text{ I/O for fetching}}$$

iii) Unclustered.

- 4 dif. values. Each semester will have approx.  $\frac{10,000}{4} = 2500$  tuples.

In the worst case

$$\text{Cost} = \boxed{\text{Some I/O for index} + 2500 \text{ I/O for fetching tuples from different pages}}$$

## 2B)

i) Table scan.

$$B(\text{Prof}) = \boxed{200 \text{ I/O}}$$

ii) Clustered hash

- There are 100 dept. In 'CENG' dept approx 10 professor.  $\left(\frac{1000 \text{ prof}}{100 \text{ dept}}\right)$

- 1000 prof  $\rightarrow$  200 pages

10 prof  $\rightarrow$  2 pages

$$\text{Cost} = \boxed{\text{Some I/O for finding correct bucket} + 1 \text{ I/O for chain, since 5 prof in the 1st bucket, the remaining 5 prof. in the second bucket}}$$

iii) Unclustered hash

- In the bucket, there will be 10 values pointing different pages.

$$\text{Cost} = \boxed{\text{Some I/O for finding correct bucket} + 10 \text{ I/O}}$$



2c)

i) Table scan

$$B(\text{prof}) = \boxed{200 \text{ I/O}}$$

ii) Clustered B+ tree

1000pid  $\rightarrow$  1000prof  $\rightarrow$  200pages

200pid  $\rightarrow$  40pages

$$\text{Cost} = \boxed{\text{Some I/O for finding the starting index (200)} + 40 \text{ I/O for linear search}}$$

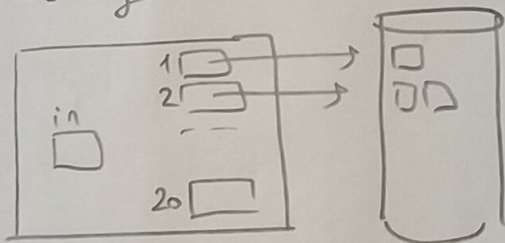
iii) Unclustered B+ tree

200pid  $\rightarrow$  200prof  $\rightarrow$  in the index leaf we will have 200 dif. values

$$\text{Cost} = \boxed{\text{Some I/O for finding the index} + 200 \text{ I/O}}$$

2d)

i) Hashing



$$\text{Cost} = \boxed{200 \text{ I/O}} \text{ for reading all professor pages}$$

+  
each partition will have approx. 5 pages, since  $5 \text{ page} \times 10 \text{ name per page} \times 20 \text{ part} = 1000 \text{ rows}$   
writing partitions to disk =  $5 \times 20 = \boxed{100 \text{ I/O}}$

$$+ \text{ read partitions and do distinct operation} = \boxed{100 \text{ I/O}}$$

$$= \boxed{400 \text{ I/O}}$$



- 2d) (Continued)

ii) Sorting

Read original relation =  $\boxed{200 \text{ I/O}}$

+

remaining relation will =  $\boxed{100 \text{ I/O}}$

have 100pgs  
(only names)

read from disk and  
sort, it will be in

one pass since  $\frac{100}{20} < 20$  ✓

=  $\boxed{100 \text{ I/O}}$

=  $\boxed{400 \text{ I/O}}$



Q3)

a)  $B(R) + B(R) \cdot B(S)$

$1000 + 1000 \cdot 1500 =$

$$\begin{array}{r} 1500,000 \\ + 1,000 \\ \hline 1,501,000 \end{array}$$

b)  $B(R) + \frac{B(R)}{(M-2)} B(S) = 1000 + \frac{1000}{100} \cdot 1500 =$

$$\begin{array}{r} 15,000 \\ + 1,000 \\ \hline 16,000 \end{array}$$

c)  $B(S) + \frac{B(S)}{(M-2)} B(R) = 1500 + \frac{1500}{100} \cdot 1000 =$

$$\begin{array}{r} 15,000 \\ + 1,500 \\ \hline 16,500 \end{array}$$

d)

i) Inner (S) has clustered index

R  $\bowtie$  S

$\begin{pmatrix} b_1 \\ b_2 \\ \vdots \end{pmatrix}$



$\begin{matrix} c_1 & b_1 \\ c_2 & b_1 \\ c_1 & b_2 \\ c_3 & b_3 \end{matrix}$

for each element

find matching

$$B(R) + T(R) \cdot \frac{B(S)}{V(S,b)}$$

$= 1000 + 10,000 \cdot \frac{1,500}{150}$

$= 101,000$

$\begin{pmatrix} 10 \end{pmatrix}$  → only need to find 10 match (in pages) ✓

ii)  $B(R) + T(R) \cdot \frac{T(S)}{V(S,b)}$

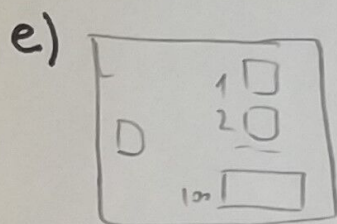
$= 1000 + 10,000 \cdot \frac{6000}{150}$

40 → since it is not clustered join attribute could be in more difo pages.

$= 401,000$



### Q3) (continued)



Each bucket will have  $\frac{B(R)}{M} = \frac{1000}{100} = 10$  pages

Since  $\frac{B(R)}{M} < M$ , each bucket will fit in the memory.

Partitioning

In R: Each bucket will have

$$\frac{B(R)}{M} = 10 \text{ pages} \quad (100 \text{ buckets} = 1000 \text{ pages})$$

In S: Each bucket will have

$$\frac{B(S)}{M} = 15 \text{ pages} \quad (100 \text{ buckets} = 1500 \text{ pages})$$

$$(\text{Read} + \text{Write}) R + (\text{Read} + \text{Write}) S = 2000 + 3000 = \boxed{5000}$$

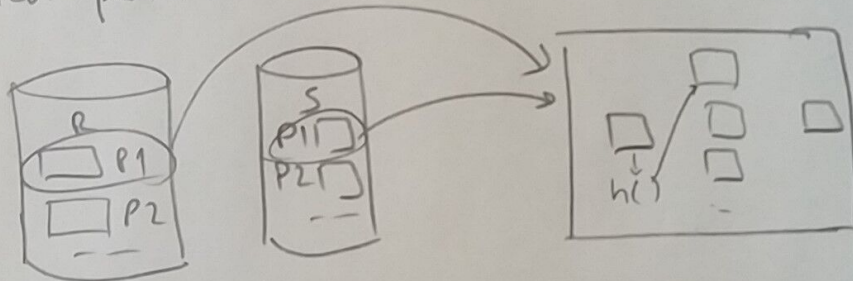
Joining  
In the final step R's and S's partitions will be read from the

disk

$$5000 + 1000 + 1500 = \boxed{7500}$$

Only needed partitions will be read.

(i.e.)



f) 1<sup>st</sup>, we have external sorting

Sorting

- Generate runs of R  $\rightarrow$  10 runs of R, last run of R is 91 pages long.

-  $S \rightarrow$  15 runs of S, last run of S is 85 pages long.

$$\text{Total cost until now} = 2B(R) + 2B(S) = \boxed{5000} \quad (1 \text{ write} + 1 \text{ read})$$

Merge-Join step

Since we will read runs of R and S

$$5000 + 1000 + 1500 = \boxed{7500}$$

$$\left. \begin{array}{l} B(R) < M^2 \\ B(S) < M^2 \end{array} \right\} \rightarrow \text{Both will fit.}$$

