Koray Can
Yurtseven
2099547

**1A)**

$$\sigma_{c>10}$$

$$\gamma_{quantity,\ Count(\delta cid)\ c}$$

$$\Pi_{quantity,\ cid}$$

$$\bowtie_{prd=prd}$$

$$\bowtie_{cid=cid}$$

$$\sigma_{prd=100}$$

Product

$$\sigma_{email\ like\ "\%gmail.com"}$$

Customer

Order

## 1B)

$\Pi_{name, phone}$

$\bowtie_{cid=crd}$

Customer     (T4)

$\Pi_{cid}$     (T3) $\Pi_{cid}$

Customer     (T2) $\bowtie_{prd=prd}$

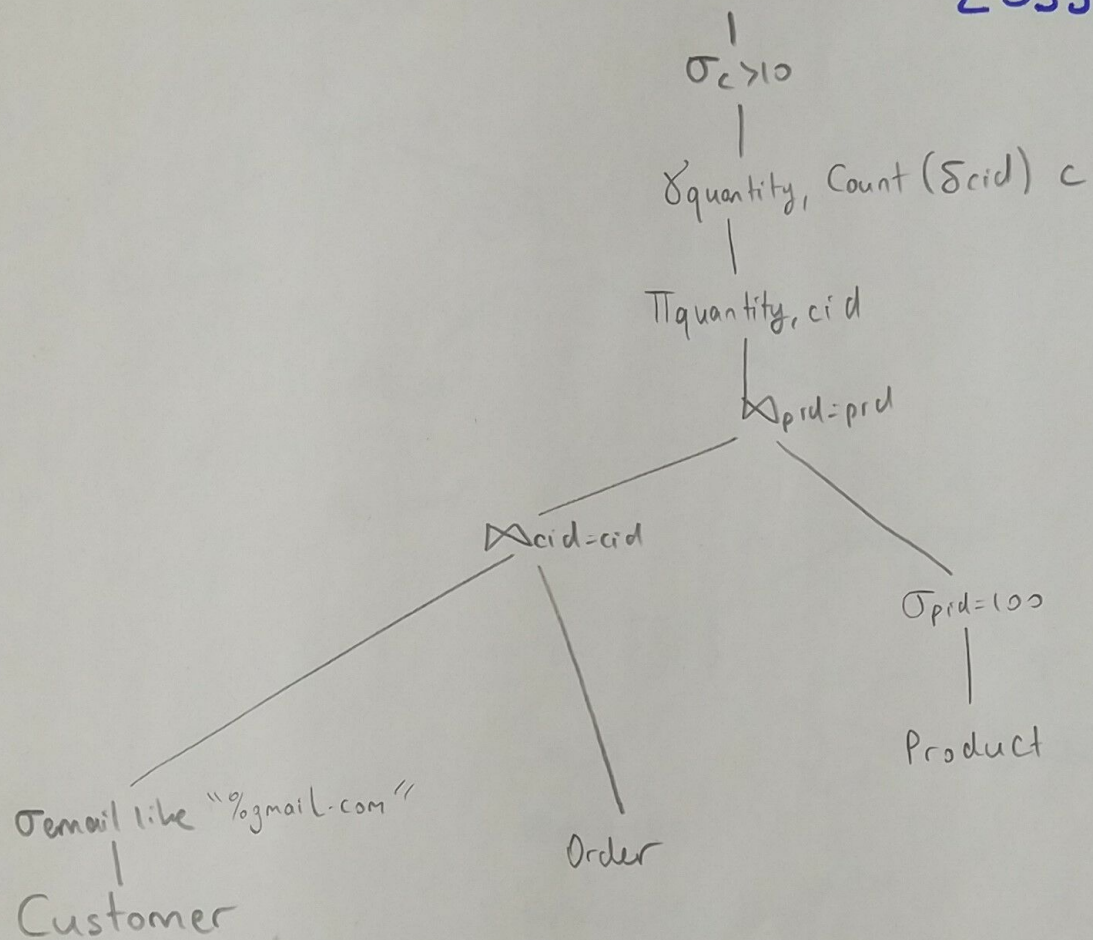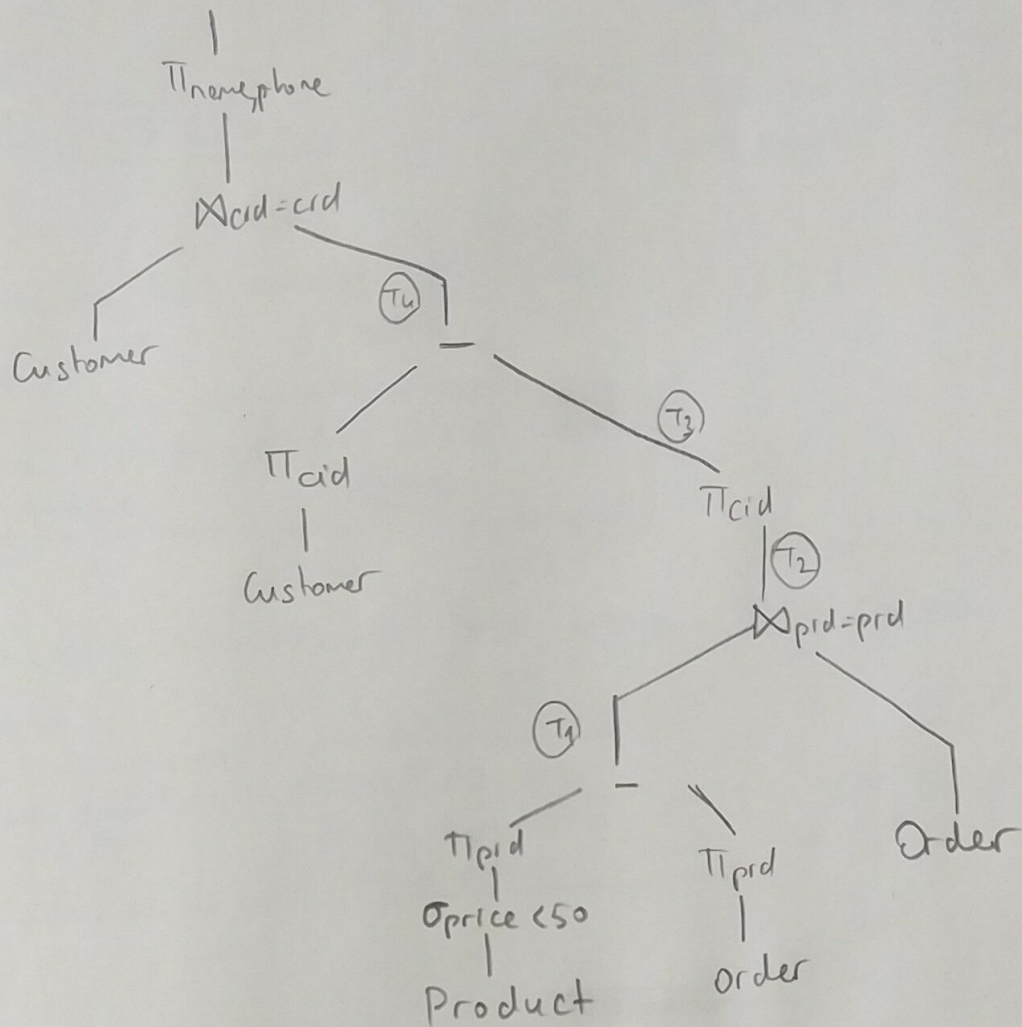(T1) $\Pi_{prd}$    $\Pi_{prd}$    Order

$\sigma_{price < 50}$    Order

Product

### Clarification:

$T_1 \Rightarrow$ $\left(\begin{array}{l}\text{Select prd} \\ \text{From Product} \\ \text{wher price} < 50 \\ \text{EXCEPT} \\ \text{Select Prd} \\ \text{From Order}\end{array}\right)$ → unnesting the inner most query → , to get cid, we need $\Rightarrow T_2$

$T_1 \bowtie$ Order

↓

Then, project Rid in $T_2 \to T_3$

outer query ←

$\left(\begin{array}{l}\text{Select cid} \\ \text{From Customer} \\ \text{Except} \\ \text{Select cid} \\ \text{From } T_3 \left(\begin{array}{l}\text{or } T_3, \text{they are} \\ \text{the same}\end{array}\right)\end{array}\right) = T_4$

In the final step

$\left(\begin{array}{l}\text{Select name, phone} \\ \text{From Customer, } T_4 \\ \text{Where } T_4.cid = \text{Customer.cid}\end{array}\right)$

②

## 2A)

i) Table scan

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

ii) Clustered B+ on semester

-4 different values. each semester will have

approx 125 page.

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

iii) Unclustered.

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

In the worst case

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

*(right margin, boxed):*
$B(Prof) = 200$
$T(Prof) = 1,000$
$B(Teach) = 500$
$T(Teach) = 10,000$

## 2B)

i) Table scan.

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

ii) Clustered hash

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

- 1000 prof → 200 pages
  10 prof → 2 pages

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

iii) Unclustered hash

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

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

③

## 2c)

### i) Table scan

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

### ii) Clustered B+ tree

$1000 pid \longrightarrow 1000 prof \longrightarrow 200 pages$

$200 prd \longrightarrow 40 pages$

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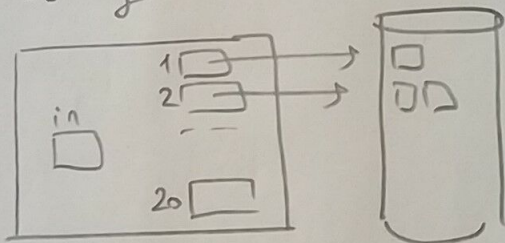
### iii) Unclustered B+ tree

$200 pid \longrightarrow 200 prof \longrightarrow$ in the index leaf we will have 200 diff. values

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

## 2d)

### i) Hashing



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

+

each partition will have approx. 5 pages, since    5 page x 10 row per x 20 partit = 1000 rows
                                                                                        page
writing partitions to disk $= 5 \times 20 = \boxed{100 \; I/O}$

+

read partitions and do distinct peration $= \boxed{100 \; I/O}$

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

- ## 2d) (Continued)

### ii) Sorting

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

$+$

remaining relation will $= \boxed{100 \text{ I/o}}$
  have    100 page
  (only names)
read from disk and
sort, it will be in
one pass since $\frac{100}{20} < 20$ ✓      $= \boxed{100 \text{ I/o}}$

$= \boxed{\underline{400 \text{ I/o}}}$

## Q3)

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

$1000 + 1000 \cdot 1500 = \boxed{1,501,000}$

$$\frac{1,500,000}{+\ 1,000}{1,501,000}$$

**b)** $B(R) + \dfrac{B(R)}{(M-2)} B(S) = 1000 + \dfrac{1000}{100} \cdot 1500 = \boxed{16,000}$

$$\frac{15,000}{+\ 1,000}{16,000}$$

**c)** $B(S) + \dfrac{B(S)}{(M-2)} B(R) = 1500 + \dfrac{1500}{100} \cdot 1000 = \boxed{16,500}$

**d)**

**i)** Inner (S) has clustered index

$R \bowtie S$



for each element     find matching

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

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

$\qquad\qquad\qquad \boxed{10} \rightarrow$ only need to find 10 match ✓ (in pages)

$= \boxed{101,000}$

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

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

$\qquad\qquad\qquad 40 \rightarrow$ since it is not clustered join attribute could be in more dif. pages.

$= \boxed{401,000}$

## e)



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

__In S__: Each bucket will have $\frac{B(S)}{M} = 15$ pages ($100$ bucket $= 1500$ pages)

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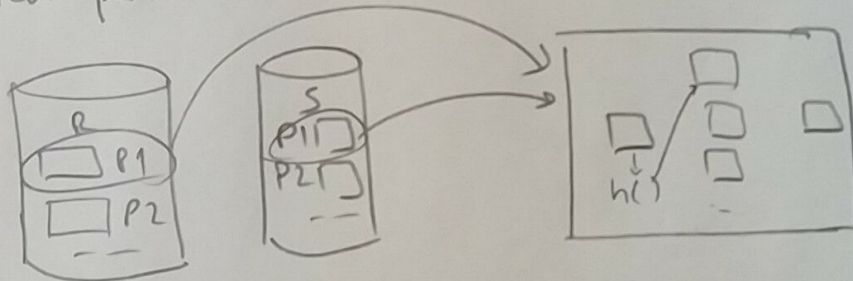
**Joining**

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) 1st, we have external sorting

**Sorting**

- Generate runs of R → 10 runs of R, last run of R is 91 pages long.
- // S → 15 runs of S, last run of S is 85 pages long.

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

**Merge-Join step**



10 for R
15 for S

Since we will read runs of R and S

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

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

⑦