

Advance Database Concepts (CS4064)

Assignment 3 Total Marks: 150

TA: Muhammad Zayam Amjad

Hasan Yahya

Name

22L-7971

Roll No.

BSE-6A

Section

Do not write below this line.

Note: Please ensure that you attempt all questions and their respective parts in the given order.

Q. No 1: Consider the following part of library database schema and the query in SQL and RA:

Book (BookID, Title, Category, Publisher, PublishYear)

Author (AuthorID, AuthorName, Gender, Email, OriginCity)

BookAuthor (BookID, AuthorID)

SELECT Title, AuthorName, Publisher
FROM Author A JOIN BookAuthor BA ON A.AuthorID=BA.AuthorID JOIN Book B ON B.BookID=BA.BookID
WHERE Gender= 'Male' AND Category= 'Education';

π Title, AuthorName, Publisher (σ Gender= 'Male' ^ Category= 'Education' (Author * BookAuthor * Book))

Your task is to optimize this query and draw the best possible query tree for this query. Take appropriate database statistics to support your answer. (Show all steps) [10]

Initial Query Tree

π Title, AuthorName, Publisher

|

δ Gender= 'Male' ^ category= 'Education'

|

\bowtie BookAuthor, BookID = Book, BookID

\bowtie Author, AuthorID = BookAuthor, AuthorID

Book

Author

BookAuthor

Heuristics

- ① Apply select early
- ② Apply project early
- ③ Combine (joins)
- ④ Replace cartesian product with join

push down

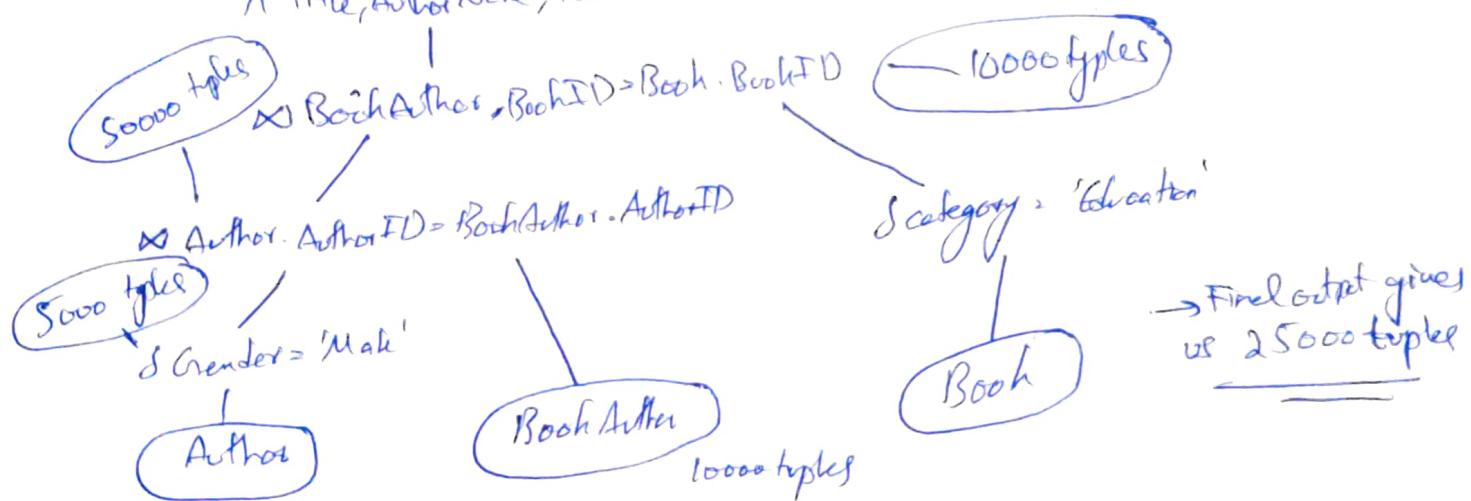
Assumptions

- Author = 10000 rows
- ↳ Gender = male 50% of them = 5000 rows
- Book = 5000 rows
- ↳ Category 'Education' is 10% of them so 500 rows
- BookAuthor = 100,000 rows

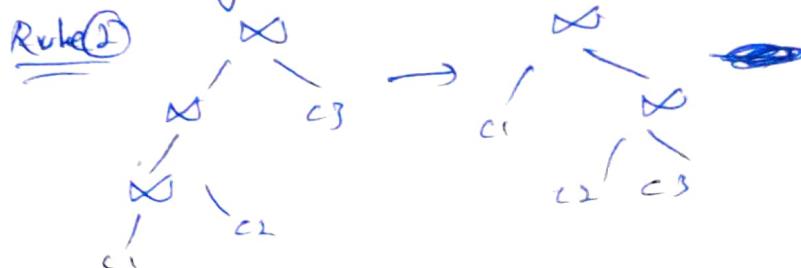
Applying Heuristics to Optimize

- ① Apply selection early

π Title, AuthorName, Publisher

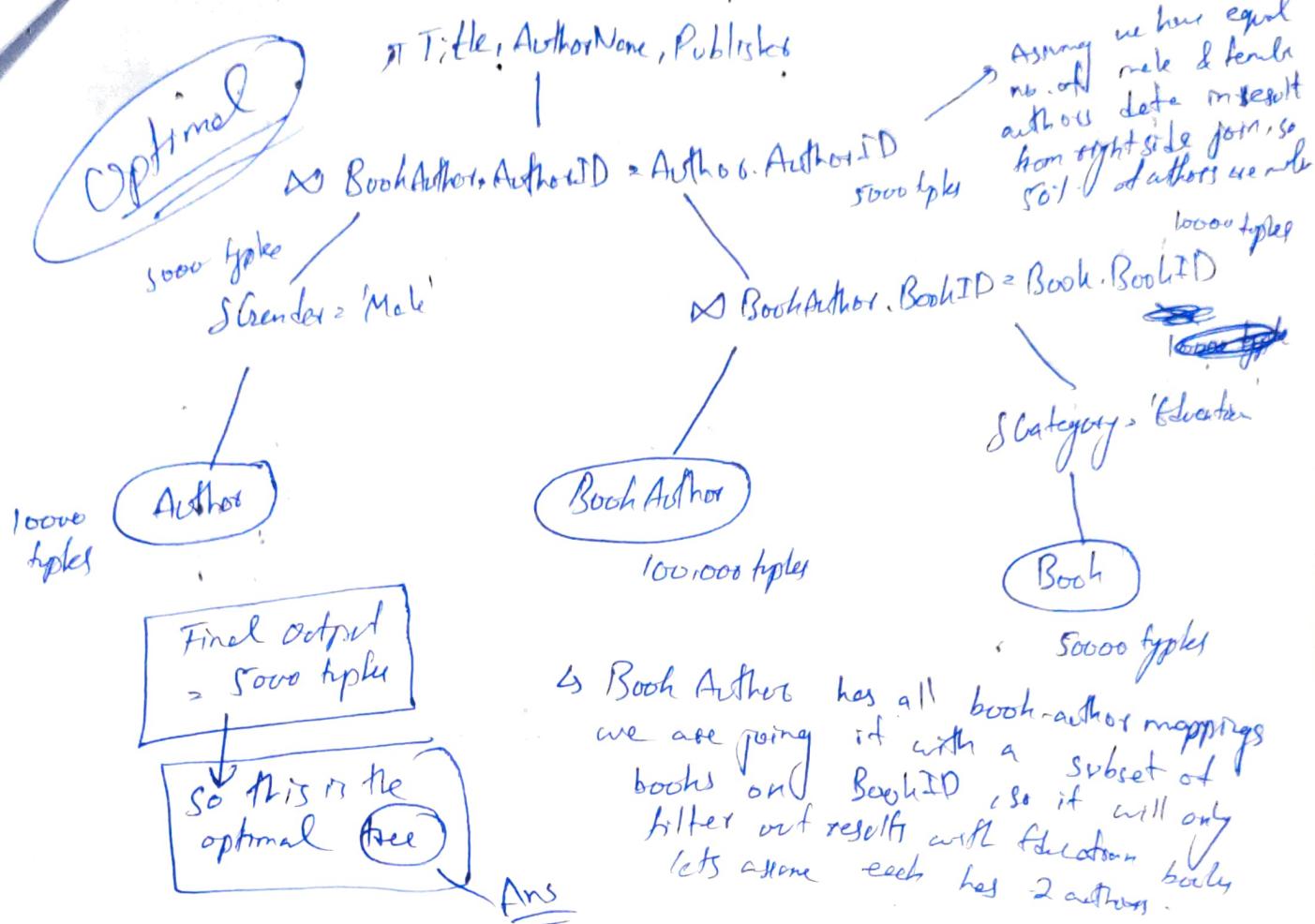


→ If we apply rules of equivalence to find possible query trees;



Rule 3 of equivalences applies as well

- $\frac{S \text{ Gender} = 'Male' \text{ (Author)}}{S \text{ Category} = 'Education' \text{ (Book)}}$



π Title, AuthorName, Publisher, σ BookAuthor. σ Category = 'Education' (Book) σ Gender = 'Male' (Author)

more optimal as category \rightarrow 'Education' is more selective than 'Gender = 'Male'' on Author and reduces number of tuples as output. Here, σ Category = 'Education' (Book) is a more selective filter.

Alternative Strategies

Next
Page

Tables	Rows	Selectivity (Filters)	Estimate Result	Result
Author	10000	Gender = 'Male' (80%)	8000	X
Book	200000	Category = 'Educational' (10%)	20000	X
BookAuthor	40,000	N/A	40000	X

→ $8000 \times 40000 \rightarrow$ small result join will give 2000 rows instead of all (80000)

(If you want to show projections at each split you can also write the same query tree as (Cophandle))

Π Title, AuthorName, Publisher



Π BookAuthor, AuthorID = Author, AuthorID

Π Title, AuthorID, Publisher



Π BookAuthor, Book = Book, BookID

Π BookID, Title, Publisher

& Category = 'Educational'

Book

Π Gender, AuthorID, AuthorName



& Gender = 'Male'

Author

(Same with extra projection specified)

CLO # 3: To develop a solution for given scenario/challenging problem in the domain of DB systems.

Q. No 2: [5+5]

- a. Consider the above library database schema and the query in SQL/RA given in Q#1. Assume that the frequency of access of this query is very high. Which attributes are more appropriate to create indexes for an efficient execution strategy to improve the performance of this query?

we want to index columns that are either in the WHERE clause or JOIN clause/conditions (to match rows efficiently). So we have 3 suggested attributes
 → Author.Gender (used in WHERE clause)
 → Book.Category (used in WHERE clause)
 → Joining ~~to~~ columns for all 3 tables (ie Author.AuthorID, Book.Author.AuthorID, Book.BookID, BookAuthor.BookID).

- b. Consider the above library database schema and assume that Book, Author and BookAuthor tables have 100000, 100000 and 50000 rows respectively. Estimate the potential Join Cardinality (jc) of $\text{Book} \bowtie_{\text{BookID}=\text{BookID}} \text{BookAuthor}$ (i.e., max number of rows resulting from the inner join of these two tables). Justify your answer.

Since there is an inner join being used and BookID is key for Book table, joining will yield at most 50000 tuples. (ie min of the 3).

As we know, join selectivity is,

$$js = 1 / \text{MAX}(\text{NDV}(\text{bookID}, \text{Book}), \text{NDV}(\text{bookID}, \text{bookAuthor}))$$

$$= 1 / 100000 = \boxed{0.00001} \text{ Ans}$$

and join cardinality is,

$$jc = js \times |\text{Book}| \times |\text{BookAuthor}|$$

$$jc = (0.00001) \times 100000 \times 50000$$

$$\boxed{jc = 50000} \text{ Ans}$$

CLO # 3: To develop a solution for given scenario/challenging problem in the domain of DB systems.

Q. No 3: Assume: A block size is $B = 1024$ bytes, file has $r = 1,000,000$ records, each record is 100 bytes long, a block pointer is $P = 10$ bytes, a record pointer is $P_R = 11$ bytes, and a key field for the index is 6 bytes long. A database system uses a B+-trees index on a key field. A leaf node and non-leaf node are one block in size and contain as many keys (and appropriate pointers) as will fit in a block. How many blocks will this index use? Also estimate the number of block accesses needed to search for and retrieve a record from the file given its key value using the B+-tree index. Show your working. [10]

Q3) Order of P :

$$(p \times P) + ((p-1) \times V_{SSN}) \leq B$$

$$(p \times 10) + ((p-1) \times 6) \leq 1024$$

$$10p + 6p - 6 \leq 1024$$

$$16p \leq 1030$$

$$p \leq 1030/16 \approx 64.375$$

So, $P = 64$

Q4) Order of P_{leaf} :

$$(P_{leaf} \times (V_{SSN} + P_R)) + P \leq B$$

$$(P_{leaf} \times (6 + 11)) + 10 \leq 1024$$

$$17P_{leaf} \leq 1024 - 10$$

$$17P_{leaf} \leq 1014$$

$$P_{leaf} \leq 1014/17$$

$$P_{leaf} \leq 59.64$$

$P_{leaf} = 59$

At level-01 (b1):

$$b1 = \left\lceil \frac{r}{p_{leaf}} \right\rceil = \left\lceil \frac{1000000}{50} \right\rceil = 16950 \text{ blocks}$$

At level-02 (b2):

$$b2 = \left\lceil \frac{16950}{64} \right\rceil = 265 \text{ blocks}$$

At level-03 (b3):

$$b3 = \left\lceil \frac{265}{64} \right\rceil = \lceil 4.14 \rceil = 5 \text{ blocks}$$

At level-04 (b4):

$$b4 = \left\lceil \frac{5}{64} \right\rceil = 1 \text{ block}$$

$$\begin{aligned} \text{Total block accesses} &= 16950 + 265 + 5 + 1 \\ &= \boxed{17221 \text{ blocks}} \end{aligned}$$

Ans

Block access cost to search a key

$$\text{value} \geq x+1 = h+1 = \boxed{5 \text{ blocks}}$$

Ans

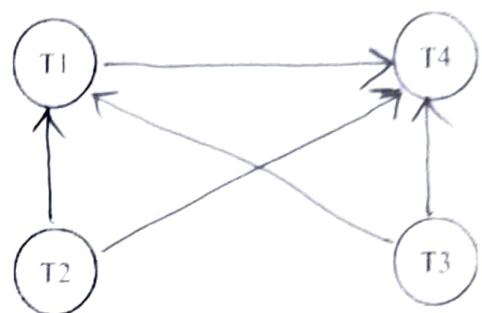
CLO # 2: Apply the models and approaches to become enabled to select and apply appropriate methods for a particular case.

Q. No 4: Consider the following schedule: [5]

S: r1(X), r2(Z), w1(Z), r3(X), r3(Y), w1(X), w3(Y), r4(Y), w4(Z), w4(Y).

Draw the serializability (precedence) graph for this schedule. State whether this schedule is conflict-serializable (correct) or not. If the schedule is conflict-serializable, write down the equivalent serial schedule(s) otherwise explain why it is not. Also state whether this schedule is view-serializable or not.

T1	T2	T3	T4
R(X)			
	R(Z)		
w(Z)		R(X)	
		R(Y)	
w(X)		w(Y)	
			R(Y)
			w(Z)
			w(Y)



Since there is no loop in the graph, it is conflict-serializable with equivalent serial schedules as $T_2 \rightarrow T_3 \rightarrow T_1 \rightarrow T_4$ and $T_3 \rightarrow T_2 \rightarrow T_1 \rightarrow T_4$.

Since all conflict serializable schedules are also view serializable, this graph is also view serializable.

CLO # 2: Apply the models and approaches to become enabled to select and apply appropriate methods for a particular case.

Q. No 5: Consider the following schedule of actions: [10+10]

S: r1(X), r2(Z), w1(Z), r3(X), r3(Y), w1(X), w3(Y), r2(Y), c3, c2, c1, w4(Z), w4(Y), c4.

For each of the following concurrency control mechanisms, describe how the concurrency control mechanism handles the schedule. Assume that the timestamp of transaction T_i is i . For lock-based concurrency control mechanisms, add lock and unlock requests to the above schedule of actions as per the locking protocol. The DBMS processes actions in the order shown. If a transaction is blocked, assume that all its actions are queued until it is resumed; the DBMS continues with the next action (according to the listed schedule) of an unblocked transaction.

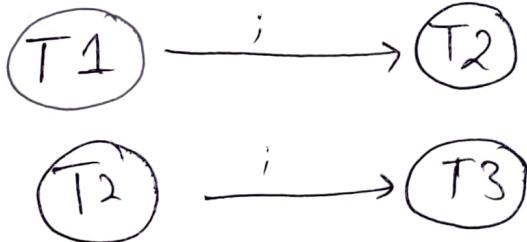
- Rigorous 2PL with timestamps used for deadlock detection (Use wait-for-graph to deal with deadlock)
- Basic Timestamp Ordering (Assume $T_1 < T_2 < T_3$)

a)

→ [S (ie s_1, s_2, \dots are for shared locks)
 x (ie x_1, x_2, \dots are for exclusive locks)]

T_1	T_2	T_3	T_4
$s_1\text{-lock}(x)$ $r_1(x)$ $x_1\text{-lock}(2)$ wait for T_2	$s_2\text{-lock}(2)$ $r_2(2)$	$s_3\text{-lock}(x)$ $r_3(x)$ $s_3\text{-lock}(y)$ $r_3(y)$ $x_3\text{-lock}(y)$ upgrade-lock ↳ $w_3(y)$ $s_2\text{-lock}(y)$ ↳ wait for T_3 ↳ $r_2(y)$ ↳ wake up restart ↳ $c_2\text{-release locks}$ $unlock_2(2)$ $unlock_2(y)$	$c_3\text{-release locks}$ $unlock_3(x)$ $unlock_3(y)$

Wait-for-graph:



$x_4\text{-lock}(2)$
$w_4(2)$
$x_4\text{-lock}(4)$
$w_4(4)$
$c_4\text{-release locks all}$
$unlock_4(2)$
$unlock_4(4)$

b)

T1	T2	T3	T4	X	Y	Z			
				RTS	WTS	RTS	WTS	RTS	WTS
R1(X)				TO	{	TO	{	TO	{
	R2(Y)				TS(T1)				
w1(Z)	→ abort T1 as RTS(2) > TS(T1)								{T2}
		R3(X)			{T3}				
		R3(Y)				{T3}			
		w3(Y)					T3		
		R2(Y)	→ abort WTS(Y) > TS(T2)						

c3

w4(Z)

w4(Y)

c4

Th

T4

with new timestamp restart T1

Restart T2

with new timestamp

imp

Note, the dots are
drawn to show which
level each entry is in

CLO # 2: Apply the models and approaches to become enabled to select and apply appropriate methods for a particular case.

Q. No 6: Consider the following log at the point of system crash. Suppose that we use ARIES recovery algorithm to answer the following questions. [9]

LSN	Last_LSN	Trans_ID	Type	Page_ID	Other_Info
1	0	T1	Update	A	...
2	0	T2	Update	C	...
3	1	T1	Commit		...
4	2	T2	Update	A	...
5	begin checkpoint				
6	end checkpoint				
7	4	T2	Commit		...
8	0	T3	Update	B	...
9	0	T4	Update	C	...
10	8	T3	Update	A	...
11	9	T4	Commit		...

Master Log: $\underline{\underline{LSN}} = 5$

- a. Show the contents of transaction and dirty page table at the time of checkpoint. [2]

Trans_ID	LSN	Status
T1	3	commit
T2	4	in progress
-	-	-

Page_ID	LSN
A	1
C	2
-	-

- b. What is done during Analysis? Be precise about the points at which Analysis begins and ends and show the contents of transaction and dirty page table reconstructed in this phase. [3]

5 (lie begin checkpoint)

Analysis Phase: From: LSN = 5 To: LSN = 11

Transaction Table

Trans_ID	LSN	Status
T1	3	commit
T2	7	commit
T3	10	in progress
T4	11	commit

Page_ID	LSN
A	5 (lie 1)
B	8
C	2

Dirty Page Table

Analysis phase starts from last checkpoint and updates the table & dirty page table which values of Transaction in fun will later be used in the coming UNDO and REDO phases.

c. What is done during Redo? Be precise about the points at which Redo begins and ends. [2]

Begin-LSN = 1

END-LSN = 10

LSNs (1, 2, 4, 8, 9, 10) update the corresponding pages (A, C, A, B, C, A). Redo does all operations from start to end except for uncommitted ones.

d. What is done during Undo? Be precise about the points at which Undo begins and ends. [2]

Begin-LSN = 10

End-LSN = 8

LSNs (10, 8) and corresponding pages (A, B), undo will only perform on T3. Undo starts from end & rollbacks uncommitted/incomplete transactions.

Q. No 7: Consider the bank database, and the following SQL query: [5+5+5]

Customer (custID, custName, cnic, birthDate, address, ...)

Account (accNo, custID, accTitle, accType, openingDate, ...)

Transaction (tID, accNo, transType, amount, transDate, ...)

Applying Deletions because it is already optimized.

Write an efficient relational-algebra expression that is equivalent to these queries and draw the optimal query plan for this query.

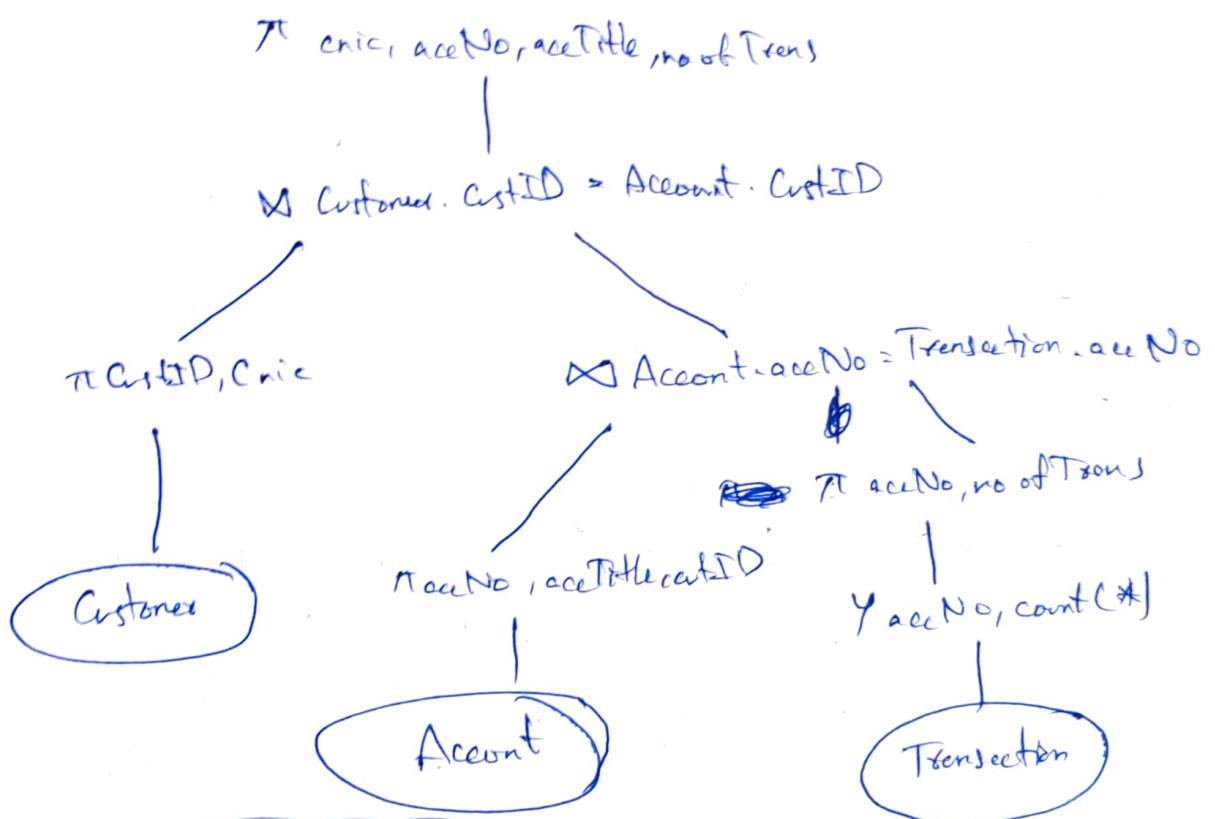
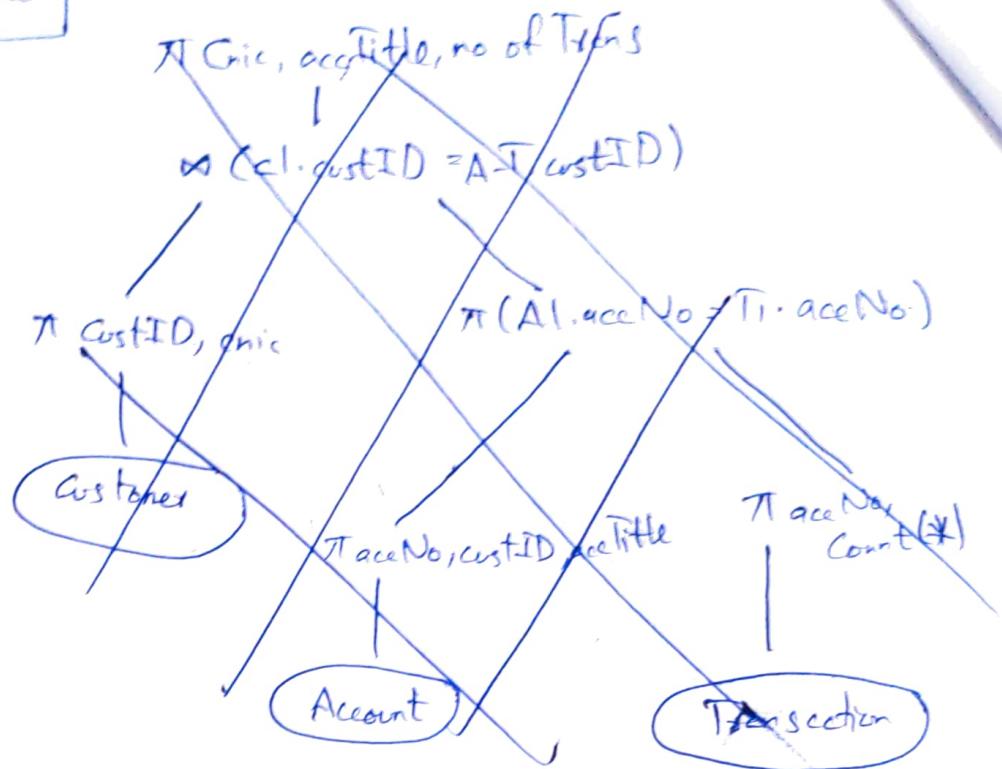
1. $\text{SELECT C.cnic, A.accNo, A.Title, T.noOfTrans FROM customer C JOIN account A ON C.custID=A.custID JOIN (SELECT accNo, COUNT(*) AS noOfTrans FROM transaction GROUP BY accNo) T ON A.accNo=T.accNo}$

RA Expression: (optimized)

$\pi_{\text{Cnic, accTitle, no of Trans}}(\pi_{\text{accNo}}(\pi_{\text{accNo, custID, accTitle}}(\text{Account}) \bowtie_{\text{accNo = accNo}} \pi_{\text{accNo:Count(*)}}(\text{Transactions})) \bowtie_{\text{CustID = CustID}} \pi_{\text{Cnic (Customer)}}(Customer))$

Next
Page

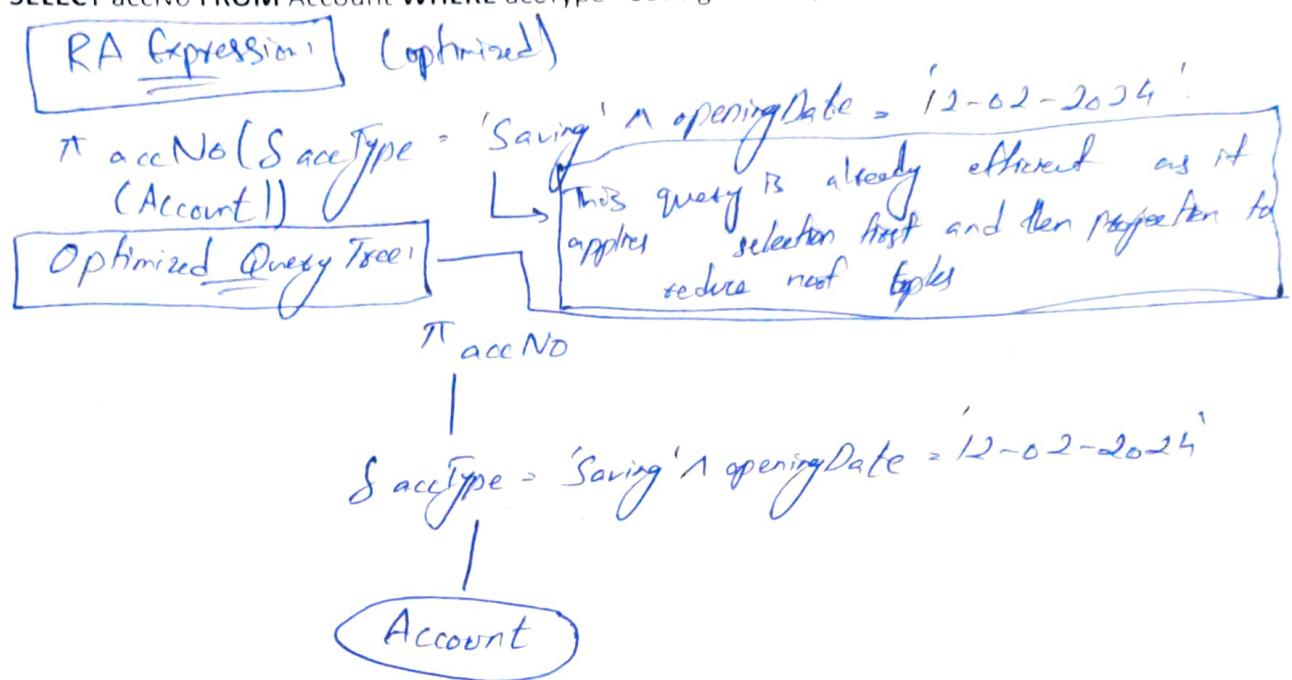
Optimized Query Tree



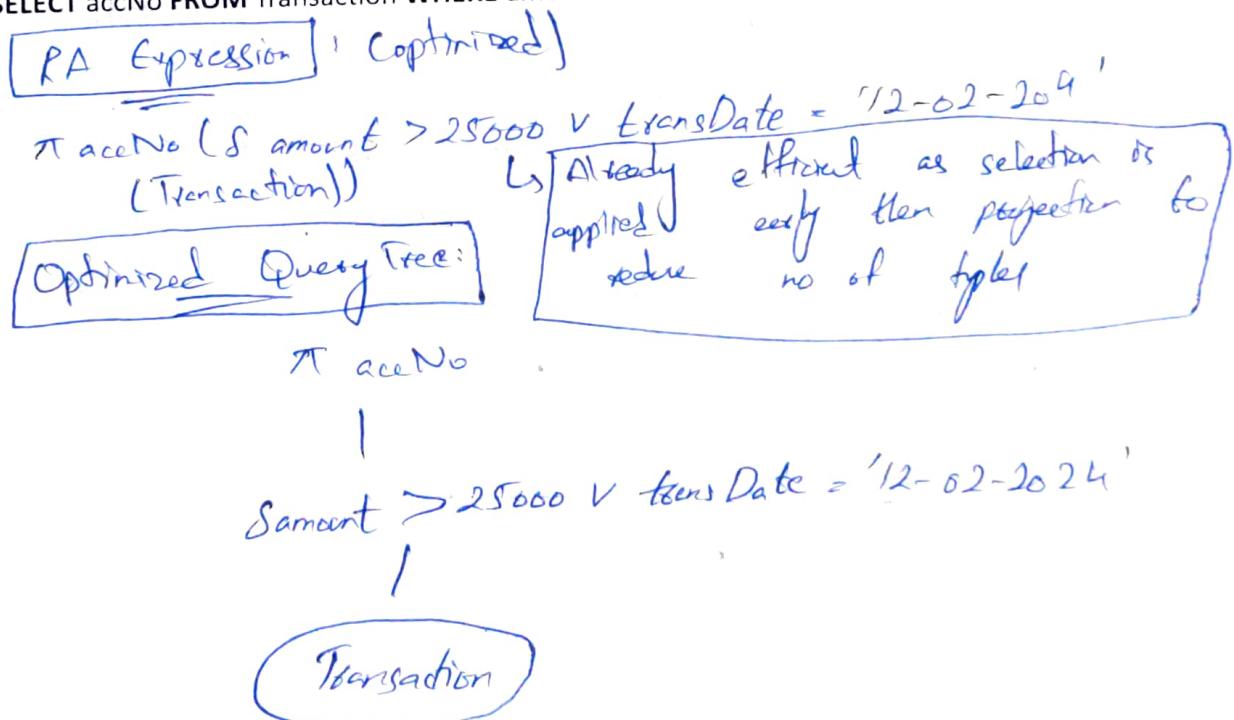
The original Query's RA tree

→ Account \bowtie (T.accNo (T.accNo, count(*) Transaction) T.accNo, count(*))
 \bowtie (accNo, count(*) Transaction))
 $\rightarrow \pi$ C.cnic, A.accNo, A.Title, T.no of Trans (Customer & Account \bowtie
 \bowtie (T.accNo (accNo, count(*) Transaction)))
 we apply Heuristics to ~~order of join~~ to get the optimal query & tree.

2. **SELECT accNo FROM Account WHERE accType='Saving' AND openingDate='12-02-2024'**



3. **SELECT accNo FROM Transaction WHERE amount > 25000 OR transDate = '12-02-2024'**



② Alternative for ②

we could trim the columns first

& then selection

& then projection.

ie
=

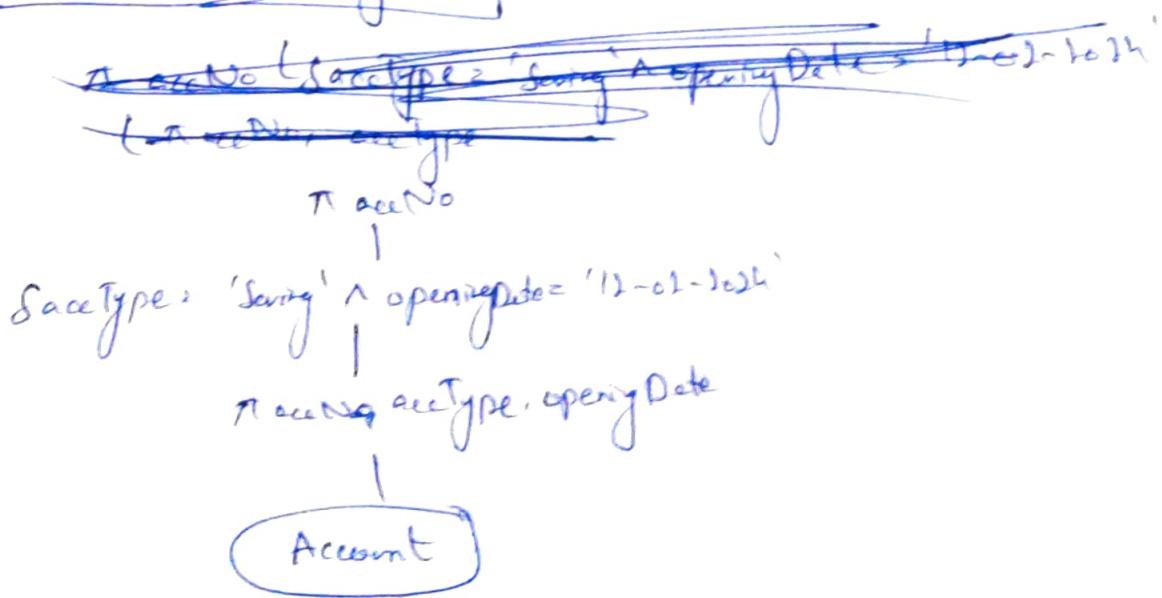
Query:

$\pi_{\text{accNo}} (\& \text{actype} = 'Savings' \wedge \text{openingDate} = '12-01-2024') (\pi_{\text{actype}},$
 $\text{openingDate} (\text{Account}))$.

This approach reduces the tuple width early on
& keeps only necessary attributes in the memory

But this strategy is only effective if Account has a lot
of columns (10-20) around. The
(---) in the Bank database shows tables
may have more columns.

Optimized Query Tree:



Q. No. 8: Assume that the number of buffers available in main memory for implementing the join is $k = n_b = 5$ blocks. The DEPARTMENT file consists of $r_D = 40$ records stored in $b_D = 10$ disk blocks and that the EMPLOYEE file consists of $r_E = 600$ records stored in $b_E = 200$ disk blocks. $[3*5=15]$

Suppose that secondary indexes exist on both the attributes

SSN of employee and Mgr_SSN of department,

with the number of index levels $X_{SSN} = 4$ and $X_{Mgr_SSN} = 2$, respectively.

$$\text{Write cost} = ((js \times |R| \times |S|) / bfrs)$$

DEPARTMENT $\bowtie_{Mgr_SSN=SSN}$ **EMPLOYEE**

Apply below Joining Algorithms and find cost estimation. Suggest what will be the best joining algorithm to apply on above case.

1. Nested-loop join (nested-block join)

if Department records in outer loop,

$$C_{j1} = b_D + (\lceil b_D / (n_B - 1) \rceil \times b_E) + ((js \times |R| \times |S|) / bfrs)$$

~~$+ 10 + (10/(5-1) \times 200) + ((4 \times 10 \times 600) / 10)$~~

$$C_{j1} = 10 + (\lceil 10 / (5-1) \rceil \times 200) + \text{write cost}$$

$$= 10 + (4 \times 200) = 10 + 800 = 810 + \text{write cost}$$

or if employee records in outer loop,

$$C_{j1} = 200 + (\lceil 200 / 3 \rceil \times 10)$$

\uparrow write cost

$\boxed{810}$

\uparrow Ans

2. Index-based nested-loop join

if Department in outer loop,

$$C_{j2a} = b_D + (1B) \times (n_B + 1S)_{SSN}$$

$\boxed{870}$

\uparrow Ans

\uparrow write cost

Department in outerloop is better.

~~$+ 10 + 400 \times (1 + 10)$~~

~~$= 10 + (400 \times 6) = 1250 + \text{write cost}$~~

Next
Page

~~\uparrow if employee records in outer loop;~~

~~$= 200 + 600 \times (2 + 1 + \frac{40}{600})$~~

if Employee records in outer loop,

$$= bE + (rE \times (X_{mng_SSN} + 1)) + \text{write cost}$$

$$= 200 + (600 \times (2 + 1)) + \text{write cost}$$

$$\Rightarrow 200 + (600 \times 3) + \text{write cost}$$

$$\Rightarrow 200 + (1800) + \text{write cost}$$

$$\Rightarrow \boxed{2000} + \text{write cost}$$

Ans

if Department records in outer loop,

$$= bD + (rD \times (X_{SSN} + 1)) + \text{write cost}$$

$$\Rightarrow 10 + (40 \times (4 + 1)) + \text{write cost}$$

$$\Rightarrow 10 + (40 \times 5) + \text{write cost}$$

$$\Rightarrow \boxed{210} + \text{write cost}$$

Ans

So. Department records in outer loop is better.

3. Sort-merge join

→ if both E and D are already sorted,

$$\text{Cost} = bE + bD + \text{write cost}$$

$$= 200 + 10 = \boxed{210} \text{ Ans} + \text{write cost}$$

→ if they are not sorted,

$$\text{Cost} = bE + bD + bE \log_2 bE + bD \log_2 bD + \text{write cost}$$

$$= 200 + 10 + 200 \log_2 (200) + 10 \log_2 (10) + \text{write cost}$$

$$= \boxed{1771.9965} \text{ Ans} + \text{write cost}$$

4. Partition-hash join

$$\text{Cost} = 3 \times (bE + bD) + \text{write cost}$$

$$= 3 \times (200 + 10) + \text{write cost}$$

$$= 3 \times (210) + \text{write cost}$$

$$= \boxed{630} \text{ Ans} + \text{write cost}$$

Conclusion:

If we can control if data is sorted or unsorted in outer loop index records in which attribute should be based nested loop joins with Department outer loop. ($\text{Cost} = 210$). However, if we do not have control over them then, the worst-case cost of partition hash (630) is better than worst case cost of enjoin.

→

Hash joins

are almost always better than sorting based join algorithms, but there are cases where sorting based joins are

preferred, (such as on non-uniform ~~data~~ data
 or where data is ~~not~~ sorted on join keys
 or where result itself needs to be
 sorted).

Condition	Sub-condition	Cost
Nested Loop (Block)	Outer = Employee	870
	Outer = Department	810
Nested Loop (Index)	Outer = Employee	2000
	Outer = Department	210
Sort Merge	Sorted	210
	Unsorted	1771.985 $= 1772$
Partition Hash	XX	630

Lowest ←

Q. No. 9: Write down the Cost function of the following selection operations [1*5=5]

1. Primary index to retrieve a single record C_{S3a} (3a)

$$C_{S3a} = x + 1$$

2. A hash key to retrieve a single record C_{S3b} (3b)

$$C_{S3b} = 1$$

3. An ordering index to retrieve multiple records C_{S4} (4)

$$C_{S4} = x + \frac{b}{2}$$

4. A clustering index to retrieve multiple records C_{S5} (5)

$$C_{S5} = x + \lceil \frac{s}{bfr} \rceil$$

5. Secondary(B+-tree) index in worst case as well as for range queries C_{S6a} & C_{S6b} (6a & 6b)

$$C_{S6a} = x + 1 + s \quad (\text{worst case})$$

$$C_{S6b} = x + \frac{b \pi}{2} + \frac{s}{2} \quad (\text{for range queries})$$

Now consider following Statistics of EMPLOYEE

Assuming $4 * 8 = 32$ this is 200.

Suppose that the EMPLOYEE file has $r_E = 1000$ records stored in $b_E = 2.00$ disk blocks with blocking factor $bfr_E = 5$ records/block and the following access paths: ~~EMPLOYEE~~

1. A clustering index on Salary, with levels $x_{\text{Salary}} = 3$ and average selection cardinality $s_{\text{Salary}} = 20$.

Find a selectivity of s_{Salary} ?

$$s_{\text{Salary}} = \frac{20}{1000} = \boxed{0.02} \quad \underline{\underline{Ans}}$$

2. A secondary index on the key attribute Ssn, with $x_{Ssn} = 4$ and average selection cardinality $s_{Ssn} = 1$, what would be the $sl_{Ssn=?}$

$$sl_{Ssn} = \frac{1}{1000} = \boxed{0.001} - \underline{\underline{Ans}}$$

$\hookrightarrow \boxed{= \frac{S_{Ssn}}{RE}}$

3. A secondary index on the nonkey attribute Dno, with $x_{Dno} = 2$ and first-level index blocks $b_{1Dno} = 4$. There are $NDV(Dno, EMPLOYEE) = 125$ distinct values for Dno, so what would be the the selectivity of Dno is $sl_{Dno} = ?$ and the selection cardinality is $s_{Dno} = ?$

$$sl_{Dno} = (1/NDV(Dno, Employee)) = \frac{1}{125} = \boxed{0.008}$$

$\underline{\underline{Ans}}$

$$S_{Dno} = (RE \times sl_{Dno}) = (RE / NDV(Dno, Employee))$$

$$= (1000 \times 0.008) = \boxed{8} - \underline{\underline{Ans}}$$

4. A secondary index on Sex, with $x_{Sex} = 1$. There are $NDV(Sex, EMPLOYEE) = 2$ values for the Sex attribute, so the average selection cardinality is $s_{Sex} = ?$

$$S_{Sex} = (RE / NDV(Sex, Employee)) = \boxed{500}$$

$\underline{\underline{Ans}}$

$$So,$$

$$S_{Sex} = (1000 / 2) = \boxed{500} - \underline{\underline{Ans}}$$

$\underline{\underline{Ans}}$

National University of Computer and Emerging Sciences
Lahore Campus

The use of cost functions with the following examples and suggest the best cost function for each operation.

OP1: $\sigma \text{SSN} = '123456789'$ (EMPLOYEE)

On this we can either use Linear Search (S1) or Secondary B+ index (S6a) for worst case

$$C_{S1b} = \left(\frac{bE}{2} \right) = \frac{200}{2} = 100$$

Best $\leftarrow C_{S6a} = XSSN + 1 = 4 + 1 = 5$ — best cost function is B+ index (worst case)
ie $C_S = XSSN + 1$.

OP2: $\sigma Dno > 5$ (EMPLOYEE)

On this we can either use Linear Search (S1) or Secondary B+ index for range (S6b),

Best $\leftarrow C_{S1a} = bE = 200$ — operation on non-key attributes.

$$C_{S6b} = XDno + \left(\frac{b1/2Dno/2}{2} \right) + (re/2)$$

$$= 2 + \left(\frac{4}{2} \right) + \left(\frac{1000}{2} \right) = 504$$

OP3: $\sigma Dno = 5$ (EMPLOYEE)

On this we can either use method S1 or S6a, ie (Linear search vs B+ index for worst case)

$$C_{S1a} = bE = 200$$

Best $\leftarrow C_{S6a} = X_{Dno} + S_{Dno} = 2 + 8 = 10$ — Best to use

OP4: $\sigma Dno = 5$ AND SALARY > 30000 AND Sex = 'F' (EMPLOYEE)

S6a ie B+ index for worst case

We can use either Linear approach, or we can use anyone of these components of selection condition so, we need to check all 4 of them.

→ Linear Approach (S1):

~~$$C_{S1a} = bE = 200$$~~ — non key (this is the best cost function)

→ Dno indexed access path (S6a):

$$C_{S6a} = XDno + SDno = 2 + 8 = 10$$

Best

→ Cost of salary indexed access path (S5):

$$C_{S5} = X_{\text{salary}} + S/bfs = 3 + \frac{500}{5}$$

→ Here 500 assumed 1000/2 i.e. half employees qualify.

$$= 103$$



→ Cost of Sex indexed access path (S6a)
B+ worst case

$$C_{S6a} = X_{\text{sex}} + S_{\text{sex}} = 1 + 500 = 501$$

Best is S6a B+ worst case
on D_{no} — Ans (i.e. S6a)

National University of Computer and Emerging Sciences

if student in outer loop

$$bS + (\lceil \frac{bS}{nb-2} \rceil \times bD) \xrightarrow{\text{Lahore Campus}} 2000 + (\lceil \frac{2000}{17} \rceil \times 10) = 22000$$

but since this is $\lceil \frac{22000}{10} \rceil > 2000$ blocks we don't use if use the other one ie 3 blocks.

Q. No. 10: Assume that the number of buffers available in main memory for implementing the join is no min = 3 blocks. The DEPARTMENT file consists of $r(D) = 50$ records stored in $b(D) = 10$ disk blocks and the Student file consists of $r(S) = 6000$ records stored in $b(S) = 2000$ disk blocks. $[1+(3*6)=19]$

- a) How many total number of blocks fetched by using Nested loop join strategy if we performed following query [1]

Select * from Department D join Student S ON D.DID=S.DID

$$\text{Total Blocks} = bD + (\lceil \frac{bD}{nb-2} \rceil \times bS) = 10 + (\lceil \frac{10}{17} \rceil \times 2000)$$

- b) What would be Sort-merge join cost if:

1. If both student and Department are already sorted

$$\text{Cost} = bD + bS = 10 + 2000 = 2010 \quad \underline{\text{Ans}}$$

2. If both are not sorted and sorted cost using external sort-merge algorithm is $(b \log_2 b)$ b is the total number of blocks in a file

$$\begin{aligned} \text{Cost} &= bD + bS + bD \log_2 bD + bS \log_2 bS \\ &= 10 + 2000 + 10 \log_2 (10) + 2000 \log_2 (2000) \\ &\approx 23974.787 \quad \underline{\text{Ans}} \end{aligned}$$

- c) Suppose the Student table has 80 distinct values over the attribute DID and the Department has 50 distinct values over the attribute DID.

Select * from Department D join Student S ON D.DID=S.DID

1. Find join selectivity ratio of both tables

$$JS = \frac{1}{\max(V_D, V_S)} = \frac{1}{\max(80, 50)} = \frac{1}{80} = 0.0125 \quad \underline{\text{Ans}}$$

2. Find join cardinality of both tables

$$\begin{aligned} JC &= r(D) \times r(S) \times JS = 50 \times 6000 \times 0.0125 \\ &= 3750 \text{ records} \quad \underline{\text{Ans}} \end{aligned}$$

Select * from Department D where D.DID NOT IN (Select S.DID from Student S)

Anti-join

3. Find join selectivity ratio of both tables

$$\begin{aligned} JS &= 1 - \text{MIN}(1, \text{NDV}(T2.y) / \text{NDV}(T1.n)) \\ &= 1 - \text{MIN}(1, 80 / 50) \\ &\Rightarrow 1 - \text{MIN}(1, 1.6) \Rightarrow 1 - 1 = 0 \quad \underline{\text{Ans}} \end{aligned}$$

4. Find join cardinality of both tables

$$JC = |T1| \times JS = 50 \times 0 \Rightarrow 0 \text{ records} \quad \underline{\text{Ans}}$$

Good Luck