

**Department of Computer Science**  
**University of Karachi Final EXAM(BSCS-F) 17<sup>th</sup> Nov 2023**

Course supervisor: Dr. Nadeem Mahmood  
 Time: 2 Hours

Course No & TITLE: BSCS- 606 "Distributed Database System"  
 Total Marks: 70

Attempt any FOUR Questions: All questions carry equal marks. Switch off your mobile phones.

**Question 1**

- a) Discuss Distributed database reference architecture and also elaborate autonomy and heterogeneity w.r.t the model
- b) Given a relation R(K;A;B;C) (where K is the key) and the following query  
 $\text{SELECT * FROM R WHERE R.A = 10 AND R.B = 15}$ 
  - What will be the outcome of running PHF on this query?
  - Does the COM MIN algorithm produce in this case a complete and minimal predicate set? Justify your answer.

**Question 2**

- a) Describe how the following can be properly modeled in the database allocation problem.
  - i. Relationships among fragments
  - ii. Query processing
  - iii. Integrity enforcement
  - iv. Concurrency control mechanisms

Simplify the following query and transform it into an optimized operator tree using the restructuring algorithm where select and project operations are applied as soon as possible to reduce the size of intermediate relations.

```
SELECT          ENAME, PNAME      FROM          EMP, ASG, PROJ
WHERE          (DUR > 12 OR RESP = "Analyst")   AND          EMP.ENO = ASG.ENO
AND           (TITLE = "Elect. Eng."        OR          ASG.PNO < "P3")
AND           (DUR > 12 OR RESP NOT= "Analyst") AND          ASG.PNO = PROJ.PNO
```

**Question 3**

- a) Some architectural models favor the definition of a global conceptual schema, whereas others do not. What do you think? Justify your selection with detailed technical arguments.
- b) Give the query graph of the following query, in SQL, on our example database:  
 $\text{SELECT ENAME, PNAME FROM EMP, ASG, PROJ WHERE DUR > 12 AND EMP.ENO = ASG.ENO AND PROJ.PNO = ASG.PNO}$   
 and map it into an operator tree.

**Question 4**

- a. Differentiate:
  - i. LAV vs GAV
  - ii. Schema vs instance matching
  - viii. Join vs semi joins
- b. Consider three sources:

Database 1 has one relation Area(Id, Field) providing areas of specialization of employees; the Id field identifies an employee.

Database 2 has two relations, Teach(Professor, Course) and In(Course, Field); Teach indicates the courses that each professor teaches and In that specifies possible fields that a course can belong to.

Database 3 has two relations, Grant(Researcher, GrantNo) for grants given to researchers, and For(GrantNo, Field) indicating which fields the grants are for.

The objective is to build a GCS with two relations: Works(Id, Project) stating that an employee works for a particular project, and Area(Project, Field) associating projects with one or more fields.

- (a) Provide a LAV mapping between Database 1 and the GCS.
- (b) Provide a GLAV mapping between the GCS and the local schemas.
- (c) Suppose one extra relation, Funds(GrantNo, Project), is added to Database 3. Provide a GAV mapping in this case.

**Question 5**

- a. Consider the join graph of figure A and the following information: size(EMP) = 100, size(ASG) = 200, size(PROJ) = 300, size(EMP 1 ASG) = 300, and size(ASG 1 PROJ) = 200. Describe an optimal join program based on the objective function of total transmission time.
- b. Consider the join graph of the following figure A, and give a program (possibly not optimal) that reduces each relation fully by semijoins

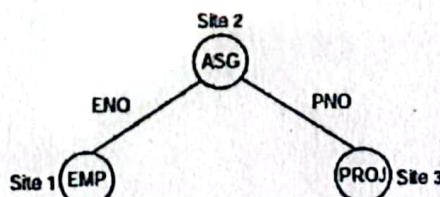
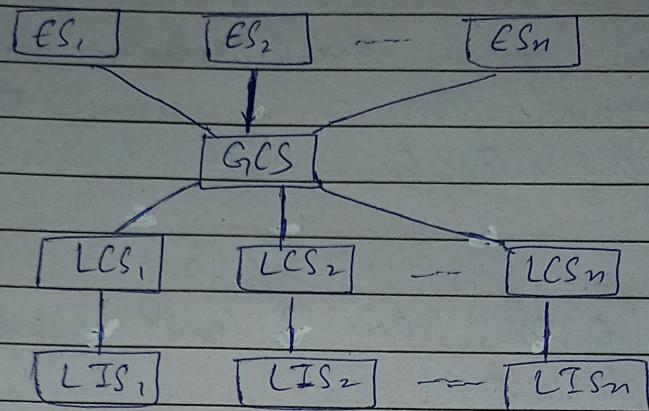


Figure A: Join Graph of Distributed Query

- Q) Discuss distributed database reference architecture & also elaborate autonomy & heterogeneity w.r.t the model.

⇒ Distributed Database Reference Architecture :-



→ ES (External Schema): The user's view of the database, shows data in a way that each user or application needs.

→ GCS (Global Conceptual Schema): Provide a unified view of all distributed data across sites.

→ LCS (Local Conceptual Schema): Describes how data is stored locally structured in each database.  
individual site.

Defines.

→ LIS (Local Internal Schema): The actual physical storage structure in each local Database

⇒ Autonomy: defines how independent each site is.

- Each site's DBMS (LCS/LIS) controls its own data & operations
- GCS coordinates across sites but does not fully control local operations.

⇒ Heterogeneity: refers to differences among sites in DBMS type, data model, or platforms.

- The architecture handles this via the global schema (GCS/ES) & communication layers, allowing integration despite local differences.

P.P(2023) Q1.b

3.10 SELECT \*

FROM R

WHERE R.A=10 AND R.B=15

(a) What will be the outcome of running PHF on this query?

PHF:  $p_1 \oplus A = 10, p_2 \oplus B = 15$

minterms  
 $m_1: p_1 \wedge p_2$      $m_2: \neg p_1 \wedge p_2$      $m_3: p_1 \wedge \neg p_2$      $m_4: \neg p_1 \wedge \neg p_2$

→ Given query will access only  $m_1: p_1 \wedge p_2 \Rightarrow A=10 \wedge B=15$  fragment.

(b) Does COM-MIN Algo produce in this case a complete & minimal predicate set?

Justify your answer.

→ Yes, COM-MIN Algo is complete (covers all data ~~in~~ query) & minimal (no unnecessary extra conditions)

→ COM-MIN returns  $\{A=10, B=15\}$  here, that set is complete & minimal for the access pattern given in the query.

**Q2.a**

Date \_\_\_\_\_

3.15: Describe how the following can be properly modeled in db allocation problem.

(a) Relationships among fragments :-

When data is fragmented, the fragments are still logically related through common attributes or keys or join constraints, to ensure the original relation can be reconstructed using joins & that referential integrity is preserved b/w related fragments.

(b) Query Processing:-

Query processing identifies where & how to execute a query when data is distributed.

It includes:

- Data localization: Identify which fragments are needed.
- Query sharding: send query to the site(s) that hold those fragments.
- Result assembly: combine partial results from different sites  
→ Goal: minimize communication cost & response time.

(c) Integrity enforcement:-

Ensures data correctness across sites. Keeps data consistent & validate across all fragments & replicated copies.

(d) Concurrency control mechanisms:-

Handle simultaneous access to data at different sites.

Prevent conflicts when many users update fragments at the same time.

→ Goal: maintain transaction consistency & avoid conflict across fragments.

→ P.P (2023)

## Q2.b

7.3 Simplify query & transform into an optimized operator tree using restructuring algo (where selection & projection operations are applied).

SELECT ENAME, PNAME

FROM EMP, ASG, PROJ

WHERE (DUR > 12 OR RESP = "Analyst")

AND EMP.ENO = ASG.ENO

AND (TITLE = "Elect. Eng")

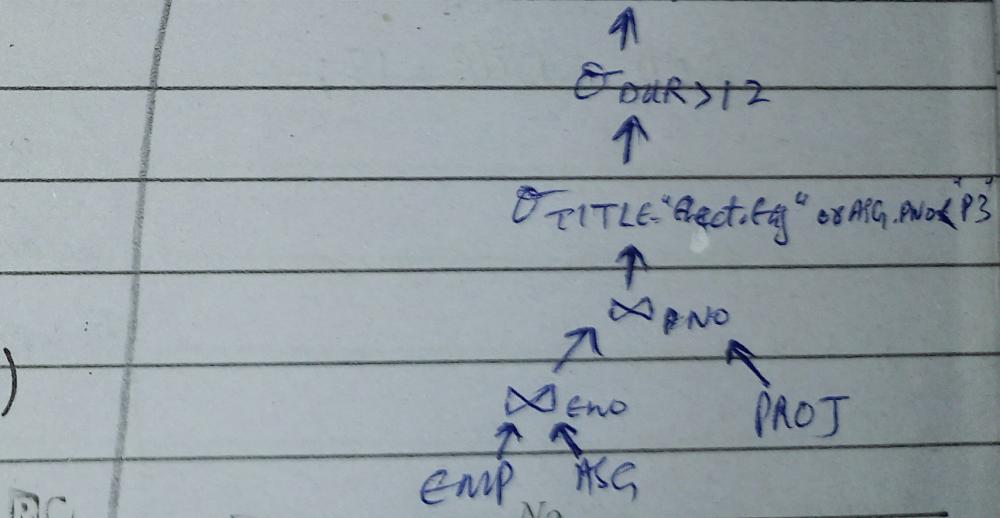
OR ASG.PNO < "P3"

AND (DUR > 12 OR RESP NOT = "Analyst")

AND ASG.PNO = PROJ.PNO;

⇒ After Simplify :-

ΠENAME, PNAME



Signature \_\_\_\_\_

RC

⇒ Generic Operator Tree.

7.3 Simplifying :-

Date \_\_\_\_\_



$\Rightarrow$  Consider only compound selection predicate in algebraic exp:-  
 $\hookrightarrow \vee, \wedge, \neg, (\text{NOT})$   $\text{EMP\_ENO} = \text{ASG\_ENO}, \text{ASG\_PNO} = \text{PROJ\_PNO}$ .

$\Rightarrow \text{Tename, PNAME} = \text{O}(\text{DUR} > 12 \vee \text{RESP} = \text{"Analyst"}) \wedge \text{EMPENO} = \text{ASG\_ENO} \wedge (\text{TITLE} = \text{"Elect. Engg"} \vee \text{ASG\_PNO} < \text{P3})$   
 $\wedge (\text{DUR} > 12 \vee \text{RESP} = \text{"Analyst"}) \wedge \text{ASG\_PNO} = \text{PROJ\_PNO}$ .

$$P_1 = \text{DUR} > 12, \quad P_2 = \text{RESP} = \text{"Analyst"}, \quad P_3 = \text{TITLE} = \text{"Elect. Engg"}, \\ P_4 = \text{ASG\_PNO} < \text{P3}.$$

$\Rightarrow$  The query qualification is:-

$$\Rightarrow (P_1 \vee P_2) \wedge (P_3 \vee P_4) \wedge (P_1 \vee \neg P_2).$$

- Acc to rules

$$\text{Rule 1: } P_1 \wedge P_2 \Leftrightarrow P_2 \wedge P_1$$

$$\Rightarrow ((P_1 \vee P_2) \wedge (P_1 \vee \neg P_2)) \wedge (P_3 \vee P_4)$$

$$\text{Rule 3: } P_1 \wedge (P_2 \vee P_3) \Leftrightarrow (P_1 \wedge P_2) \vee (P_1 \wedge P_3).$$

$$\Rightarrow (P_1 \wedge (P_2 \vee \neg P_2)) \wedge (P_3 \vee P_4)$$

$$\text{Rule 8: } P \vee \neg P \Leftrightarrow \text{true}$$

$$\text{Rule 3: } P \wedge \text{true} \Leftrightarrow P$$

$$\Rightarrow (P_1 \wedge \text{true}) \wedge (P_3 \vee P_4)$$

$$\Rightarrow \boxed{P_1 \wedge (P_3 \vee P_4)}$$

$$\text{or } (\text{DUR} > 12) \wedge (\text{TITLE} = \text{"Elect. Engg"} \vee \text{ASG\_PNO} < \text{P3}).$$

$\Rightarrow$  can't simplify further.

4.2 Some architecture models favor the definition of a GCS (global conceptual schema), whereas other do not. What do you think? Justify your selection with detailed technical arguments.

=> Global Conceptual Schema (GCS): Unified view of all data in a distributed system.

- Models favoring GCS: Centralized distributed dbs (eg: banking).
  - ↳ Ensures consistency, integrity & simpler queries.
- Models not favoring GCS: Multidatabase systems (eg: hospitals (autonomous sites))
  - ↳ sites are heterogeneous & independent; global schema is hard to maintain.

↳ In short: Use GCS when you need control & consistency; avoid it when flexibility & independence are more important.

P.P(2023)

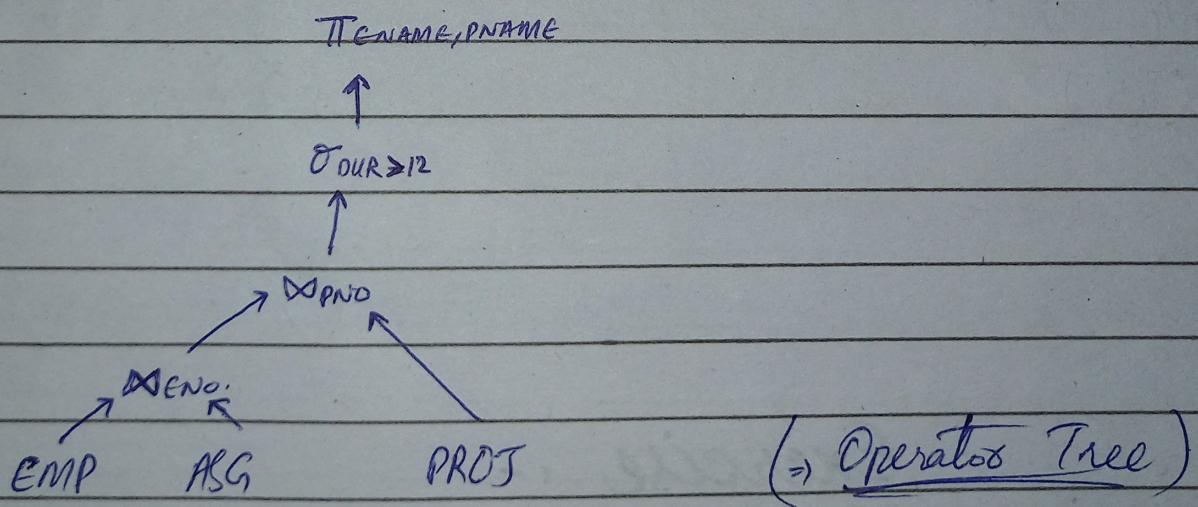
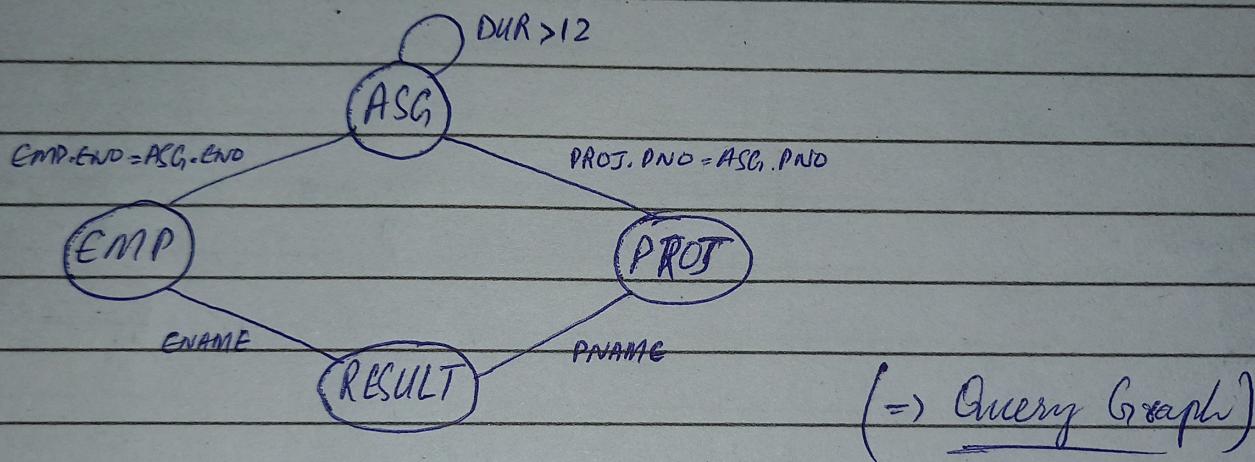
Q3.b

Date \_\_\_\_\_

7.2 - SQL Query :-

```
SELECT ENAME, PNAME  
FROM EMP, ASG, PROJ  
WHERE DUR > 12  
AND EMP.ENO = ASG.ENO  
AND PROJ.PNO = ASG.PNO
```

Give Query graph & map it into an operator tree.



→ 1.1 (2023)

## Q4.a

Date \_\_\_\_\_

Differentiate:-

1) LAV vs GAV:

- In LAV (Local-as-view): GCS already exists, each local DB are views of a global schema. (based on global).
- In GAV (Global-as-view): Build GCS by combining all local DBs. (Build from local DBs).

2) Schema vs instance matching:

- Schema Matching: matching table names or column(attribute) names. (like "Name with FullName")
- Instance Matching: matching actual data inside the tables. (like "Aliza" in one table & "Aliza Khan" in another).

3) Join vs Semi-Join:

- Join: A join combines two tables & shows data from both.

↳ Ex:

Table A	
ID	Name
1	Alize
2	Sara

Table B	
ID	City
1	Karachi
3	Lahore

query:  $A \bowtie_{A.ID = B.ID} (B)$

ID	Name	City
1	Aliza	Karachi

• Semi-Joins:

↳ Only shows rows from the 1st table that have a matching in the second table.

↳ Ex - query:  $A \bowtie_{A.ID = B.ID} (B)$  → only keeps rows from A that match B.

Table A	
ID	Name
1	Aliza

$\rightarrow$  P.P(2023)  
4.5, 4.12

## Q4.b

Given:

↳ GCS relations:

- Works (Id, Project)  $\rightarrow$  employee works on a project.
- Area (Project, Field)  $\rightarrow$  project belongs to a field.

↳ Databases:

DB1: Area (Id, field)

DB2: Teach (Professor, Course), In (Course, field).

DB3: Grant (Researcher, Grant NO), For (Grant NO, field).

(a) LAV mapping b/w DB1 & the GCS.

LAV = view over GCS.

$\Rightarrow$  DB1. Area (Id, field)  $\rightarrow$  Works (Id, Project), Area (Project, Field).

Date \_\_\_\_\_

(b) GLAV mapping b/w the GCS & Local Schemas.

GLAV = mix of LAV & GAV

→ 1st, GCS(Work(Id, Project) & Area(Project, field)) as 1 GCS.

⇒ Works(Id, Project), Area(Project, field) → GCS(Id, field).

→ Now local schemas in GCS(Id, field) mapping-

⇒ GCS(Id, field) → DB2.Area(Id, field), DB2.Tech(<sup>Professor Id</sup>~~Course~~, Course), DB2.In(Course, field),  
DB3.Grant(Researcher<sup>Id</sup>, GrantNo), DB3.Foo(GrantNo, field).

(c) GAV mapping. one extra funds(GrantNo, project), added in DB3.

GAV = view over local

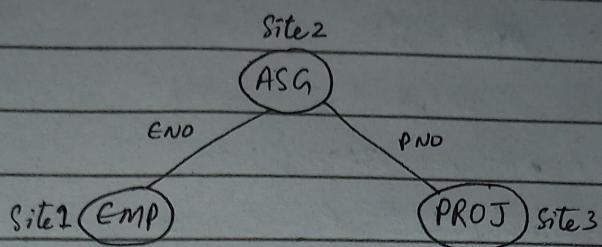
⇒ GCS.Works(Id, Project) → DB3.Grant(Researcher, GrantNo), DB3.funds(GrantNo, Project)

⇒ GCS.Area(Project, field) → DB3.Foo(GrantNo, field), DB3.funds(GrantNo, Project).

8.2.1. → P.P (2023)

### Q5.a

Consider join graph.



& following info:-

- size(EMP) = 100
- size(ASG) = 200
- size(PROJ) = 300
- size(EMP  $\bowtie$  ASG) = 300
- size(ASG  $\bowtie$  PROJ) = 200.

Based on func. of total transmission time describe an optimal join program.

⇒ Optimal Join Program (Minimizing Total transmission Time):

1) Step 1: Join ASG  $\bowtie$  PROJ at site 3. ( $ASG \rightarrow PROR$ ).  
• size( $ASG \bowtie PROJ$ ) = 200.

2) Step 2:- Join intermediate result with EMP at site 3.

- EMP  $\bowtie$  ( $ASG \bowtie PROJ$ )
- size(EMP)  $\bowtie$  size( $ASG \bowtie PROJ$ ) = 100 + 200

⇒ Total transmission Time = 300

⇒ Always join tables producing smaller intermediates first & move the ~~more~~ smaller to reduce transmission cost.

P.P(2023)

## Q5.b

Q.4. Consider 8.2 graph. & give a program that reduces each relation by semijoins.

→ Semijoin:  $R \times_{\{A\}} S =$  means keep from R, only tuples whose A-value appear in S(matching).

- Site 1  $\rightarrow$  EMP(eno)
- Site 2  $\rightarrow$  ASG (eno, pno)
- Site 3  $\rightarrow$  PROJ(pno).

→ Reduce each relation using semijoin before the final join.

• Step 1:- Reduce EMP (By removing those tuples which doesn't exist in ASG)

$$\Rightarrow EMP' = EMP \times_{eno} (ASG)$$

• Step 2 : Reduce PROJ (By removing those tuples which doesn't exist in ASG)

$$\Rightarrow PROJ' = PROJ \times_{pno} (ASG)$$

Step 3 - Reduce ASG (By using Reduce EMP & PROJ by removing those ENO's & PNO's which doesn't match with ASG(ENO, PNO)).  
Date \_\_\_\_\_

$$\Rightarrow ASG' = (ASG \times_{ENO} (EMP')) \times_{PNO} PROJ$$

Finally -

1)  $EMP' \rightarrow$  Site 2

2)  $PROJ' \rightarrow$  Site 2

3) Site 2  $\rightarrow$   $EMP' \times_{ENO} ASG' \times_{PNO} PROJ'$

$\Rightarrow$  Semijoin = reduce the size of each relation before sending or joining.  
Saves communication & computation.

$\alpha$