

Chatty \Rightarrow Query Processing

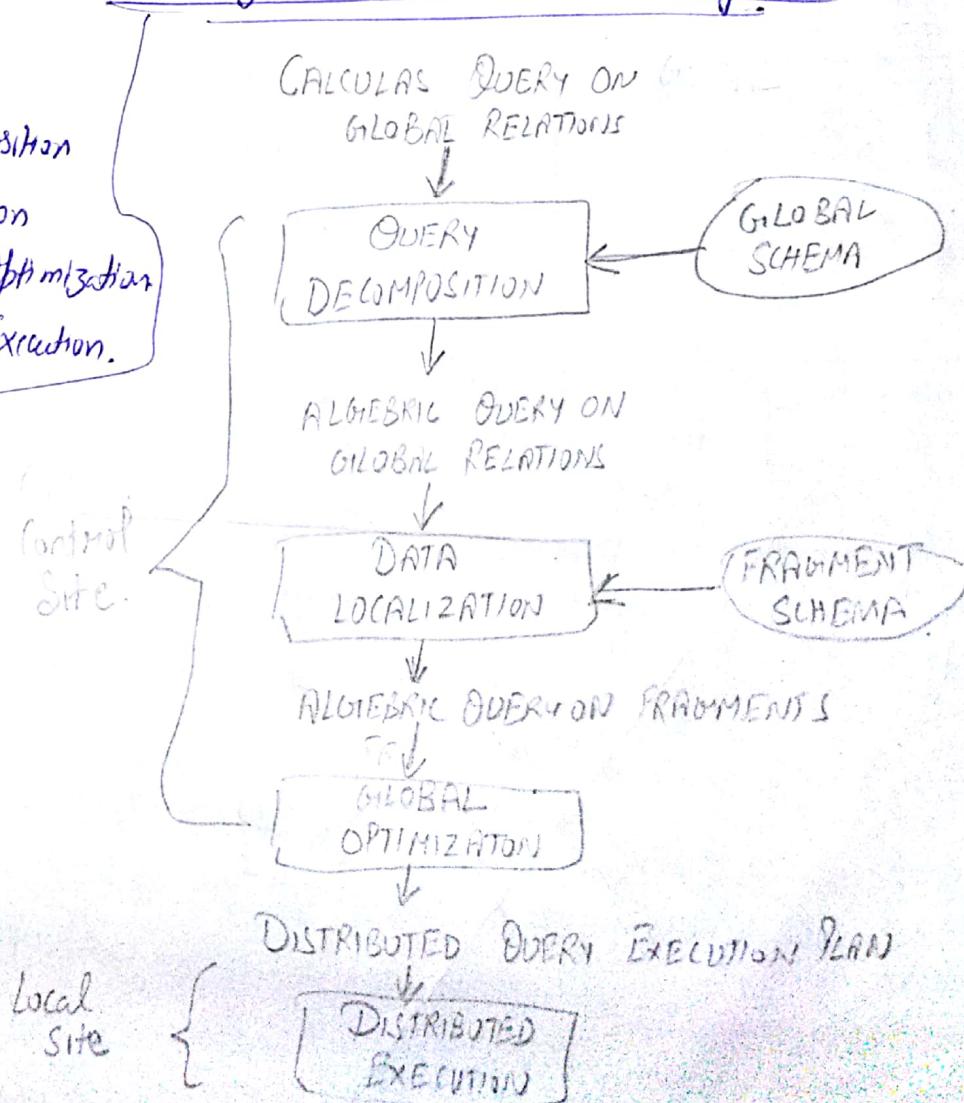
Layers of Query Processing.

1. Processing of query in distributed DBMS instead of Local (centralized) DBMS.
2. It is very difficult, because there are many (elements) factors involved in it which affect the overall performance of distributed queries.
3. That is why query processing problem is divided into several subproblems/layers, which are easier to solve individually.
4. Main function of Query Processor is to transform a high level query (calculus) into an equivalent low-level query (algebra).

Main Layers of Query Processing.

• 4 Layers

- 1) Query Decomposition
- 2) Data Localization
- 3) Global Query Optimization
- 4) Distributed Execution.



- The input is a query on global data expressed in relational calculus.
- The first three layers ~~map~~^{convert} the input query into low-level query (algebraic query) and map it into optimized distributed execution plan.
- The first three layers are performed by a central control site and use Schema, information stored in global directory.
- Fourth layer performs distributed query execution by executing the plan and return the answer to the query it is done by local sites.

Layers

1 Query Decomposition

- The first layer decomposes the calculus query into an algebraic query on ~~global~~. The information needed for this transformation is found in Global Schema.
- Query decomposition can be viewed as four successive steps:
 - Normalization → (CNF) Conjunctive Normal Form [Rewrite query into Normal Form] → (DNF) Disjunctive Normal Form.
 - Analysis → [Normalized query is analyzed ~~semantically~~]
 - Elimination of Redundancy → [Simplified the query by eliminating redundant data]
 - Rewriting → [The calculus query is reconstructed as an algebraic query. Several algebraic queries ~~can~~ can be derived from the same calculus query]

1) Data Localisation

- Output of first layer is an algebraic query on global relation which is input to the second layer.
 - The main role of this layer is to localize the query's data using data distribution information. (Fragment Scheme).
 - We know that each fragment is stored at different site.
- This layer determines which fragments are involved in the query and transforms the ~~global~~ query into fragment query.

3) Global Query Optimization

- The input of the third layer is a fragment algebraic query.
- The goal of this layer is to find an execution strategy for the algebraic fragment query which is closer to optimal.
- Query Optimization consists of:
 - 1) Finding the best ordering of operations in the fragment query.
 - 2) Find the communication operations which minimizes cost function.

4) Distributed Execution

- The last layer is performed by all the sites having fragments involved in the query.
- Each subquery called a local query, is executed at one site using local scheme.

1 Query Decomposition :

I Normalization

- The input query may be arbitrary complex
 - Transform the query to a normalized form to facilitate further processing.
- With relational languages such as SQL,
The most important transformation is that of query qualification.
(where clause)

There are two possible normal forms:

I Conjunctive NF

giving precedence to AND (\wedge)

~~Useless in
Redundancy
Removes~~
~~Redundancy~~

$$(P_1 \vee P_{1,2} \vee \dots \vee P_{1,n}) \wedge \dots \wedge (P_m \vee P_{m,2} \wedge \dots \wedge P_{mn})$$

2) Disjunctive NF

giving Precedence to (V) OR.

$$(P_1 \wedge P_{1,2} \wedge \dots \wedge P_{1,n}) \vee \dots \vee (P_m \wedge P_{m,2} \wedge \dots \wedge P_{mn})$$

In disjunctive NF, the query can be processed as independent conjunctive subqueries, linked by OR.

The transformation of predicates using well known equivalence rules for logical operations (\wedge , \vee and \neg);

- 1) $P_1 \wedge P_2 \Leftrightarrow P_2 \wedge P_1$ } Commutative law
- 2) $P_1 \vee P_2 \Leftrightarrow P_2 \vee P_1$ }
- 3) $P_1 \wedge (P_2 \wedge P_3) \Leftrightarrow (P_1 \wedge P_2) \wedge P_3$ } Associative law
Disjunction
- 4) $P_1 \vee (P_2 \vee P_3) \Leftrightarrow (P_1 \vee P_2) \vee P_3$ }
- 5) $P_1 \wedge (P_2 \vee P_3) \Leftrightarrow (P_1 \wedge P_2) \vee (P_1 \wedge P_3)$ } Distributive law
- 6) $P_1 \vee (P_2 \wedge P_3) \Leftrightarrow (P_1 \vee P_2) \wedge (P_1 \vee P_3)$ }
- 7) $\neg(P_1 \wedge P_2) \Leftrightarrow \neg P_1 \vee \neg P_2$ } Negation
- 8) $\neg(P_1 \vee P_2) \Leftrightarrow \neg P_1 \wedge \neg P_2$ }
- 9) $\neg(\neg P) \Leftrightarrow P$ }

Example : Find the names of employees who have been working on Project P_1 for 12 or 24 months.

• SOL Query \Rightarrow {

```

SELECT ENAME
FROM EMP, ASST
WHERE EMP.ENO = ASST.ENO
AND ASST.PNO = "P1"
AND DUR = 12 OR DUR = 24
    
```

}

E

ASG1, ASG2

Emp(E_{no}, E_{name}, title)

Asg(E_{no}, P_{no}, resp, d_{ur})

Proj(P_{no}, P_{name}, Budget)

CNF

The qualification in CNF +

EMP.E_{NO} = ASG1.E_{NO} \wedge ASG1.P_{NO} = "P₁" \wedge DUR = 12 \vee DUR = 24

2 DISJUNCTIVE NF (DNF)

The qualification in DNF

(EMP.E_{NO} = ASG1.E_{NO} \wedge ASG1.P_{NO} = "P₁") \wedge DUR = 12) \vee

(EMP.E_{NO} = ASG1.E_{NO} \wedge ASG1.P_{NO} = "P₁") \wedge DUR = 24)

Step
2^o

Analysis

- Query Analysis enables rejection of normalized queries for which further processing is either impossible or unnecessary.
- The main reason for rejection of query
 - ① Semantically incorrect (type incident)
 - ② Syntactically wrong (D)
- ⇒ ② A query is semantically incorrect if its component do not contribute to generation of the results.

Ex ①

Type incident means

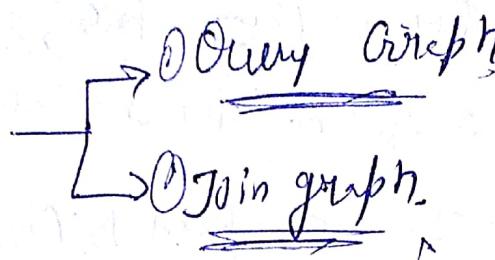
Ex ②

```
SELECT E#
  FROM EMP
 WHERE ENAME > 200
```

- ① attribute E# is not defined in scheme
- ② operation > 200 is incompatible with the type string of ENAME

Sample (2)

Semantically
incorrect



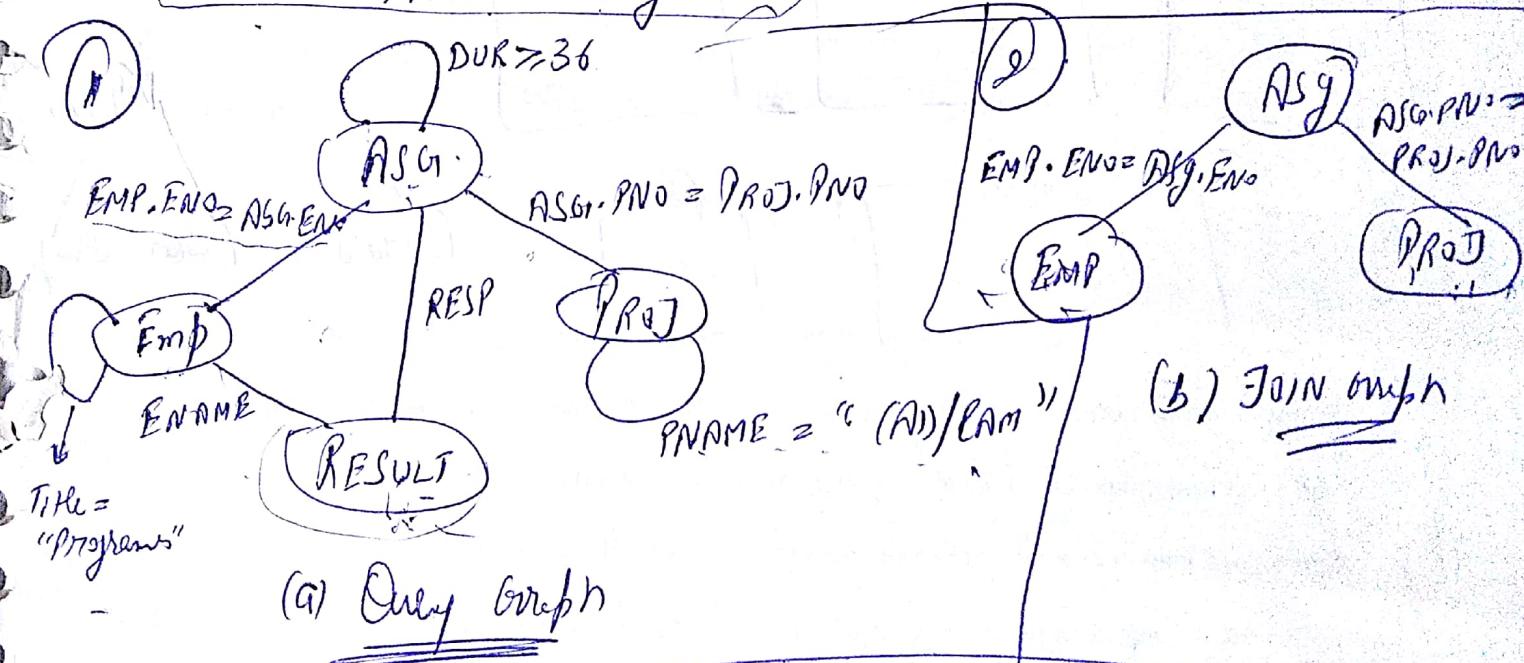
Find name, responsibility of programmers who have been working on CAD/CAM project for more than 3 years.

Example

```

SELECT ENAME, RESP
FROM EMP, ASG, PROJ
WHERE EMP.ENO = ASG.ENO
AND ASG.PNO = PROJ.PNO
AND PNAME = "CAD/CAM"
AND DUR >= 36
AND TITLE = "Programmer"
  
```

Aug(PNo, E, DUR)



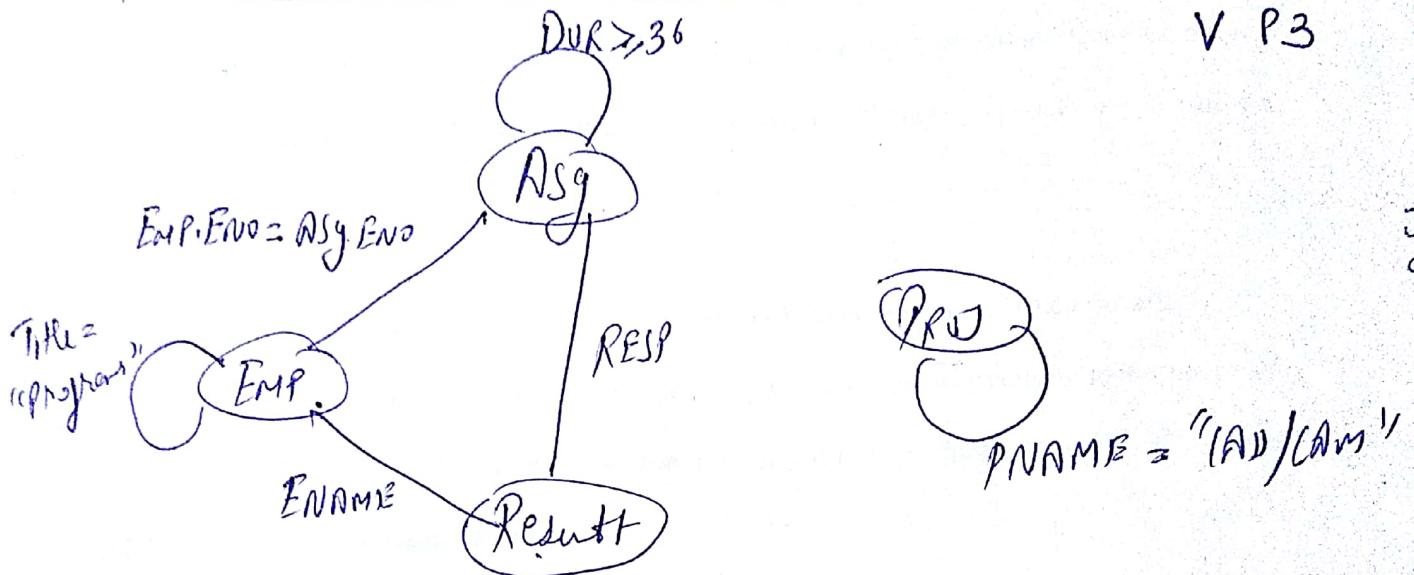
- In Query graph (connection graph) one node indicates result relation. and any other node indicates operator relation.
- Query is semantically incorrect if any node does not contribute in result
- In general join predicate is mostly ~~irrelevant~~ and the query should be rejected.

$P_1 \rightarrow P_{\text{reg}}$
 $P_{\text{reg}} \text{ else}$
 $J \cdot \text{else } P_3$
 ()
 item
 u

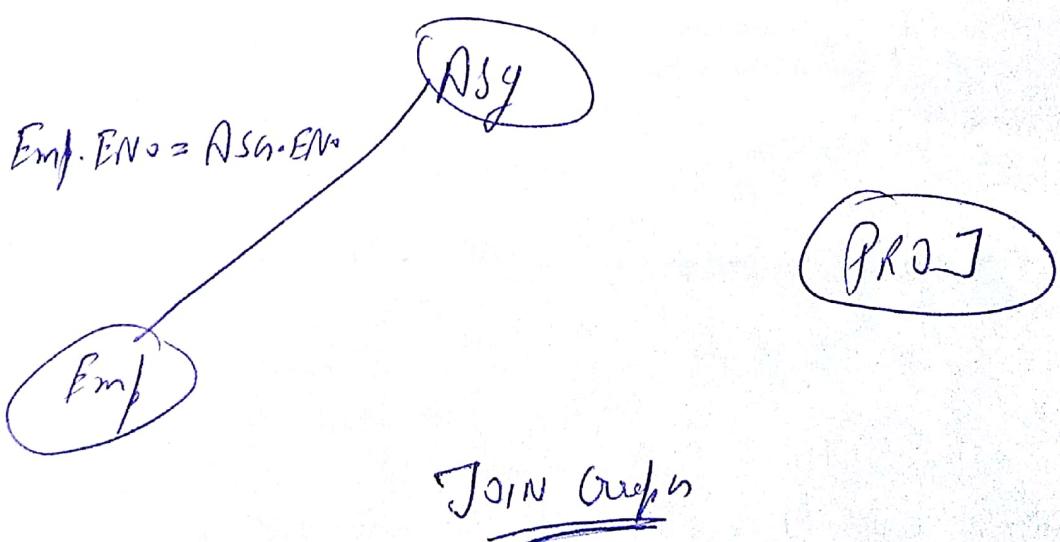
Example:
 SELECT ENAME, RESP
 FROM EMP, ASG, PROJ.
 WHERE EMP.ENO = ASG.ENO
 AND PNAME = "AD/CAM"
 AND DUR > 36
 AND TITLE = "Programmer"

$(\neg(P_1) \wedge (P_1 \vee P_2)) \wedge \neg(P_2)$,
 $\vee P_3$

Join
ca



Disjointed Query Graph
 Means Query Semantically Incorrect



Elimination of Redundancy

• Simplify the Query by eliminating Redundancy

By using following idempotency rules:-

1. $P \wedge P \Leftrightarrow P$
2. $P \vee P \Leftrightarrow P$
3. $P \wedge \text{True} \Leftrightarrow P$
4. $P \wedge \text{False} \Leftrightarrow \text{false}$
5. $P \vee \text{true} \Leftrightarrow \text{true}$.
6. $P \vee \text{false} \Leftrightarrow P$
7. $P \wedge \neg P \Leftrightarrow \text{false}$
8. $P \vee \neg P \Leftrightarrow \text{true}$
9. $P_1 \wedge (P_1 \vee P_2) \Leftrightarrow P_1$,
10. $P_1 \vee (P_1 \wedge P_2) \Leftrightarrow P_1$

Eg

```
SELECT TITLE
  FROM EMP
 WHERE NOT (TITLE = "Programmer")
   AND (TITLE = "Elect. Engg. Progammer")
   OR   TITLE = "Elect. Engne"
   AND NOT (TITLE = "Elect. Engne")
   OR   ENAME = "J. Doe"
```

can be simplified

```
SELECT TITLE
  FROM EMP
 WHERE ENAME = "J. Doe"
```

Simplification

Let,

$P_1 \Leftrightarrow \text{Title} = \text{"Programmer"}$

$P_2 \Leftrightarrow \text{Title} = \text{"Elect. Eng"}$

$P_3 \Leftrightarrow \text{Ename} = \text{"J. Doe"}$

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

Applying DeMorgan's

Disjunctive Normal Form for this Qualification is obtained by applying Rule 5 of Section 1 (Normalization)

$$\underline{\text{Step 1}} \quad (\neg P_1 \wedge ((P_1 \wedge \neg P_2) \vee (P_2 \wedge \neg P_2))) \vee P_3$$

Applying Rule 3 of Section 1

$$\underline{\text{Step 2}} \quad (\neg P_1 \wedge P_1 \wedge \neg P_2) \vee (\neg P_1 \wedge P_2 \wedge \neg P_2) \vee P_3$$

By applying rule 7 of this section,

$$\underline{\text{Step 3}} \quad (\text{false} \wedge \neg P_2) \vee (\neg P_1 \wedge \text{false}) \vee P_3$$

By applying same rule, we get

$$\underline{\text{Step 4}} \quad \text{false} \vee \text{false} \vee P_3$$

which is equivalent to P_3 by rule 4

$$\underline{\text{Step 5}} \quad P_3$$

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

Step 1 \rightarrow Rule 5

$$\boxed{P_1 \wedge (P_2 \vee P_3) \Leftrightarrow (P_1 \wedge P_2) \vee (P_1 \wedge P_3)}$$

$$\begin{array}{c} \text{Solve} \quad (\neg P_1 \wedge (P_1 \vee P_2) \wedge \neg P_2) \vee P_3 \\ \hline \end{array}$$

$$\neg P_2 \wedge (P_1 \vee P_2)$$

$$\begin{array}{c} \Downarrow \\ (\neg P_2 \wedge P_1) \vee (\neg P_2 \wedge P_2) \end{array}$$

$$\underline{(\neg P_2 \wedge P_1) \vee (\neg P_2 \wedge \neg P_2)}$$

$$\begin{array}{c} \text{Solve} \quad (\neg P_1 \wedge ((P_1 \wedge P_2) \vee (P_2 \wedge \neg P_2))) \vee P_3 \\ \hline \end{array}$$

Step 2
Rule 3

$$\underline{P_1 \wedge (P_2 \wedge P_3)} \Leftrightarrow \underline{(P_1 \wedge P_2) \wedge P_3}$$

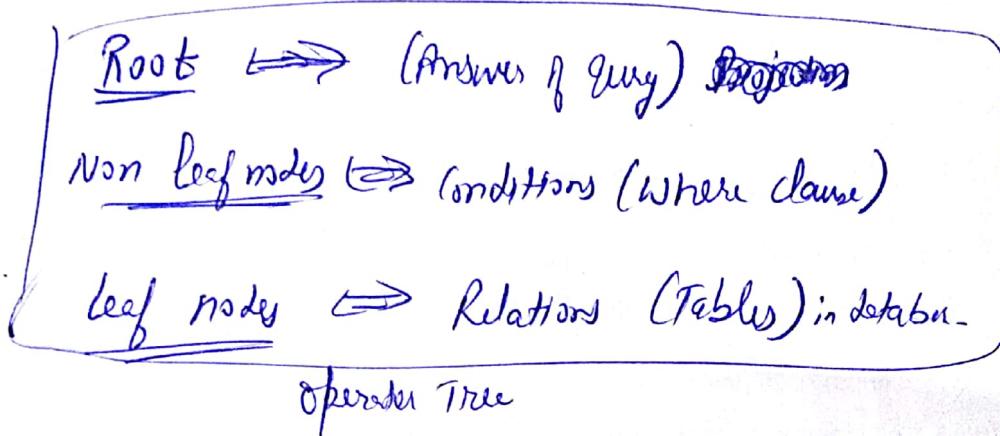
$$\underline{(\neg P_1 \wedge (P_1 \wedge P_2)) \vee \neg P_1 \wedge (P_2 \wedge \neg P_2)} \vee P_3$$

$$\begin{array}{c} \text{Solve} \quad ((\neg P_1 \wedge P_1 \wedge P_2) \vee (\neg P_1 \wedge P_2 \wedge \neg P_2)) \vee P_3 \end{array}$$

Rewriting

The last step of Query Decomposition is to rewrite the query in relational algebra.

For the sake of clarity, represent the relational algebra query graphically by operator tree.



Example

SELECT ENAME
FROM PROJ, ASG, EMP
WHERE ASG.ENO = EMP.ENO
AND ASG.PNO = PROJ.PNO
AND ENAME != "J.Doe"
AND PROJ.PNAME = "(AD/CAM)"
AND (DUR = 12 OR DUR = 24)

