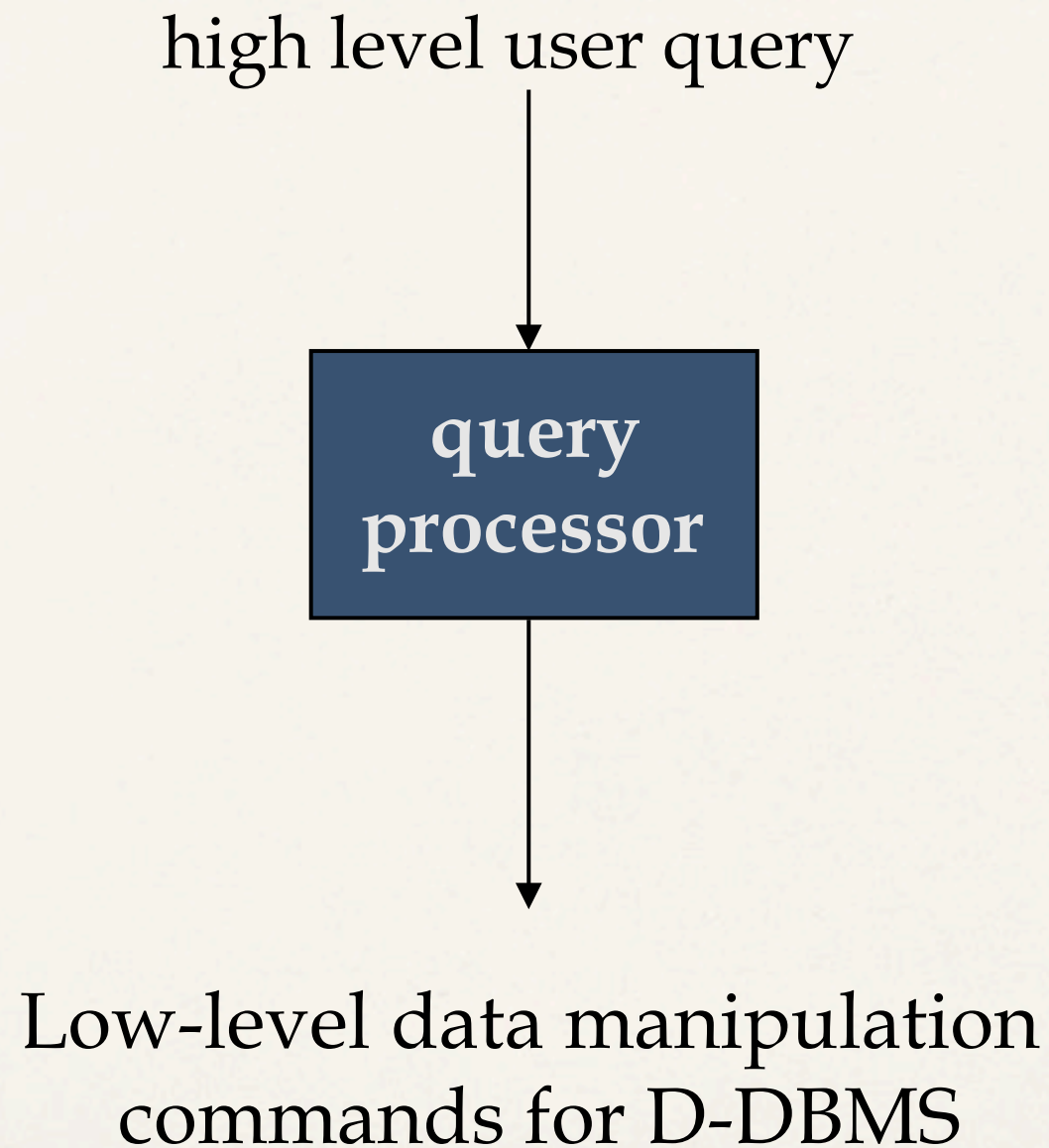


Outline

- Introduction
- Background
- Distributed Database Design
- Database Integration
- Semantic Data Control
- Distributed Query Processing
 - ➔ Overview
 - ➔ Query decomposition and localization
 - ➔ Distributed query optimization
- Multidatabase Query Processing
- Distributed Transaction Management
- Data Replication
- Parallel Database Systems
- Distributed Object DBMS
- Peer-to-Peer Data Management
- Web Data Management
- Current Issues

Query Processing in a DDBMS



Query Processing Components

- Query language that is used
 - ➔ SQL: “intergalactic dataspeak”
- Query execution methodology
 - ➔ The steps that one goes through in executing high-level (declarative) user queries.
- Query optimization
 - ➔ How do we determine the “best” execution plan?
- We assume a homogeneous D-DBMS

Selecting Alternatives

```
SELECT      ENAME
FROM        EMP, ASG
WHERE       EMP.ENO = ASG.ENO
AND        RESP = "Manager"
```

Strategy 1

$$\Pi_{ENAME}(\sigma_{RESP="Manager" \wedge EMP.ENO=ASG.ENO}(EMP \times ASG))$$

Strategy 2

$$\Pi_{ENAME}(EMP \bowtie_{ENO} (\sigma_{RESP="Manager"}(ASG)))$$

Strategy 2 avoids Cartesian product, so may be “better”

What is the Problem?

Site 1

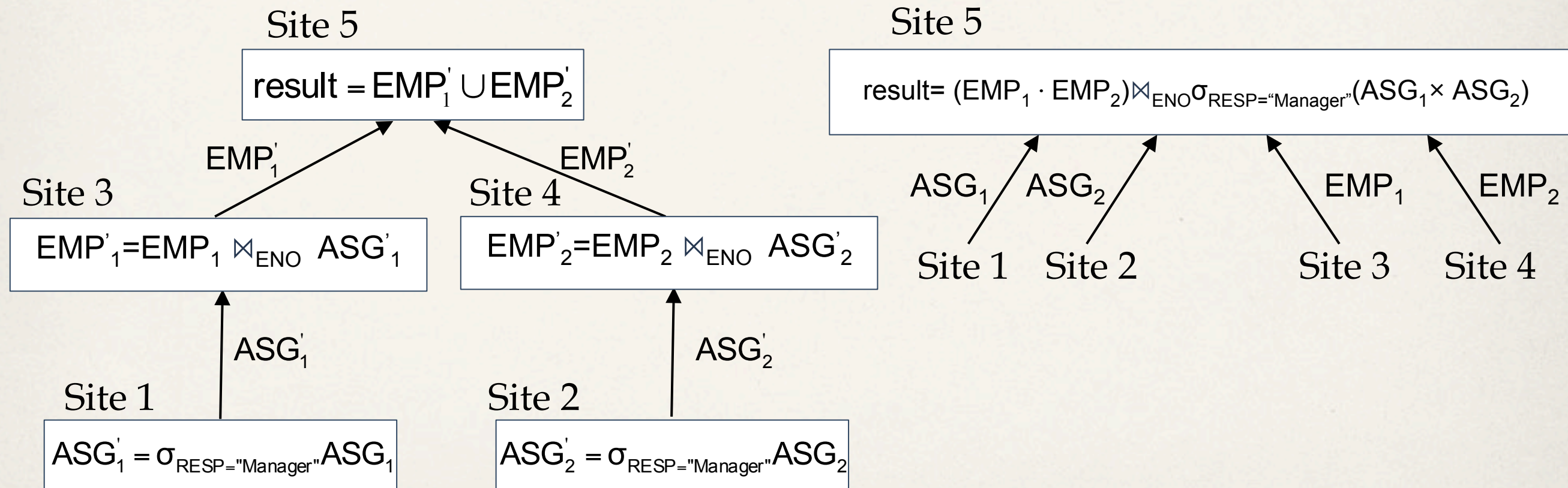
Site 2

Site 3

Site 4

Site 5

$ASG_1 = \sigma_{ENO \leq "E3"}(ASG)$ $ASG_2 = \sigma_{ENO > "E3"}(ASG)$ $EMP_1 = \sigma_{ENO \leq "E3"}(EMP)$ $EMP_2 = \sigma_{ENO > "E3"}(EMP)$ Result



Cost of Alternatives

- Assume

- $size(EMP) = 400, size(ASG) = 1000$
- tuple access cost = 1 unit; tuple transfer cost = 10 units

- Strategy 1

- produce ASG': $(10+10) * \text{tuple access cost}$ 20
- transfer ASG' to the sites of EMP: $(10+10) * \text{tuple transfer cost}$ 200
- produce EMP': $(10+10) * \text{tuple access cost} * 2$ 40
- transfer EMP' to result site: $(10+10) * \text{tuple transfer cost}$ 200
- Total Cost** 460

- Strategy 2

- transfer EMP to site 5: $400 * \text{tuple transfer cost}$ 4,000
- transfer ASG to site 5: $1000 * \text{tuple transfer cost}$ 10,000
- produce ASG': $1000 * \text{tuple access cost}$ 1,000
- join EMP and ASG': $400 * 20 * \text{tuple access cost}$ 8,000
- Total Cost** 23,000

Query Optimization Objectives

- Minimize a cost function

I/O cost + CPU cost + communication cost

These might have different weights in different distributed environments

- Wide area networks
 - ➔ communication cost may dominate or vary much
 - ◆ bandwidth
 - ◆ speed
 - ◆ high protocol overhead
- Local area networks
 - ➔ communication cost not that dominant
 - ➔ total cost function should be considered
- Can also maximize throughput

Complexity of Relational Operations

- Assume
 - relations of cardinality n
 - sequential scan

Operation	Complexity
Select Project (without duplicate elimination)	$O(n)$
Project (with duplicate elimination) Group	$O(n * \log n)$
Join Semi-join Division Set Operators	$O(n * \log n)$
Cartesian Product	$O(n^2)$

Query Optimization Issues – Types Of Optimizers

- Exhaustive search
 - ➔ Cost-based
 - ➔ Optimal
 - ➔ Combinatorial complexity in the number of relations
- Heuristics
 - ➔ Not optimal
 - ➔ Regroup common sub-expressions
 - ➔ Perform selection, projection first
 - ➔ Replace a join by a series of semijoins
 - ➔ Reorder operations to reduce intermediate relation size
 - ➔ Optimize individual operations

Query Optimization Issues – Optimization Granularity

- Single query at a time
 - ➔ Cannot use common intermediate results
- Multiple queries at a time
 - ➔ Efficient if many similar queries
 - ➔ Decision space is much larger

Query Optimization Issues – Optimization Timing

- Static
 - ➔ Compilation ➔ optimize prior to the execution
 - ➔ Difficult to estimate the size of the intermediate results \Rightarrow error propagation
 - ➔ Can amortize over many executions
 - ➔ R*
- Dynamic
 - ➔ Run time optimization
 - ➔ Exact information on the intermediate relation sizes
 - ➔ Have to reoptimize for multiple executions
 - ➔ Distributed INGRES
- Hybrid
 - ➔ Compile using a static algorithm
 - ➔ If the error in estimate sizes $>$ threshold, reoptimize at run time
 - ➔ Mermaid

Query Optimization Issues – Statistics

- Relation
 - ➔ Cardinality
 - ➔ Size of a tuple
 - ➔ Fraction of tuples participating in a join with another relation
- Attribute
 - ➔ Cardinality of domain
 - ➔ Actual number of distinct values
- Common assumptions
 - ➔ **Independence** between different attribute values
 - ➔ **Uniform distribution** of attribute values within their domain

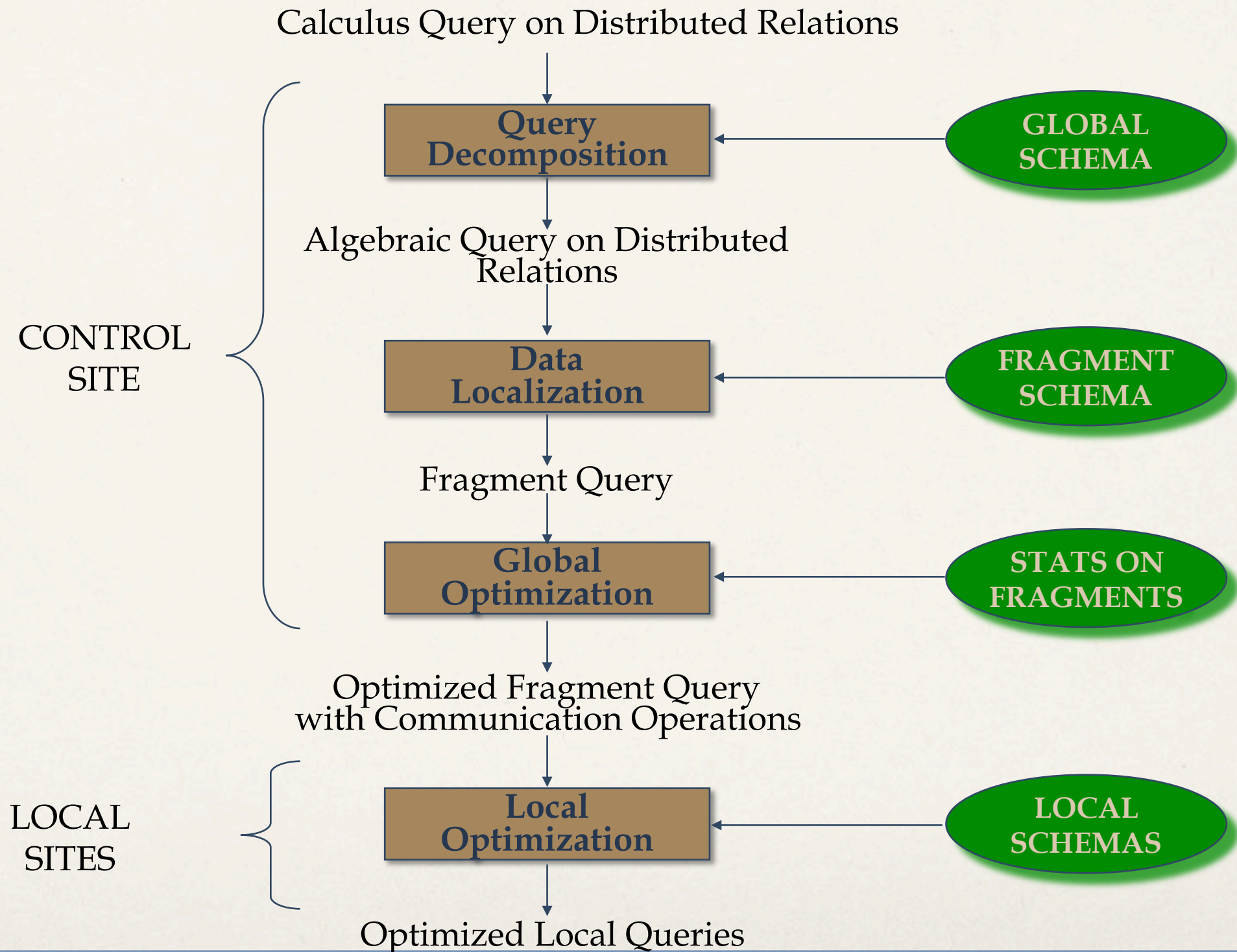
Query Optimization Issues – Decision Sites

- Centralized
 - ➔ Single site determines the “best” schedule
 - ➔ Simple
 - ➔ Need knowledge about the entire distributed database
- Distributed
 - ➔ Cooperation among sites to determine the schedule
 - ➔ Need only local information
 - ➔ Cost of cooperation
- Hybrid
 - ➔ One site determines the global schedule
 - ➔ Each site optimizes the local subqueries

Query Optimization Issues – Network Topology

- **Wide area networks** (WAN) – point-to-point
 - ➔ Characteristics
 - ◆ Low bandwidth
 - ◆ Low speed
 - ◆ High protocol overhead
 - ➔ Communication cost will dominate; ignore all other cost factors
 - ➔ Global schedule to minimize communication cost
 - ➔ Local schedules according to centralized query optimization
- **Local area networks** (LAN)
 - ➔ Communication cost not that dominant
 - ➔ Total cost function should be considered
 - ➔ Broadcasting can be exploited (joins)
 - ➔ Special algorithms exist for star networks

Distributed Query Processing Methodology



Outline

- Introduction
- Background
- Distributed Database Design
- Database Integration
- Semantic Data Control
- Distributed Query Processing
 - ➔ Overview
 - ➔ Query decomposition and localization
 - ➔ Distributed query optimization
- Multidatabase query processing
- Distributed Transaction Management
- Data Replication
- Parallel Database Systems
- Distributed Object DBMS
- Peer-to-Peer Data Management
- Web Data Management
- Current Issues

Step 1 – Query Decomposition

Input : Calculus query on global relations

- Normalization
 - ➔ manipulate query quantifiers and qualification
- Analysis
 - ➔ detect and reject “incorrect” queries
 - ➔ possible for only a subset of relational calculus
- Simplification
 - ➔ eliminate redundant predicates
- Restructuring
 - ➔ calculus query ➔ algebraic query
 - ➔ more than one translation is possible
 - ➔ use transformation rules

Normalization

- Lexical and syntactic analysis
 - ➔ check validity (similar to compilers)
 - ➔ check for attributes and relations
 - ➔ type checking on the qualification

- Put into **normal form**

- ➔ Conjunctive normal form

$$(p_{11} \vee p_{12} \vee \dots \vee p_{1n}) \wedge \dots \wedge (p_{m1} \vee p_{m2} \vee \dots \vee p_{mn})$$

- ➔ Disjunctive normal form

$$(p_{11} \wedge p_{12} \wedge \dots \wedge p_{1n}) \vee \dots \vee (p_{m1} \wedge p_{m2} \wedge \dots \wedge p_{mn})$$

- ➔ OR's mapped into union
- ➔ AND's mapped into join or selection

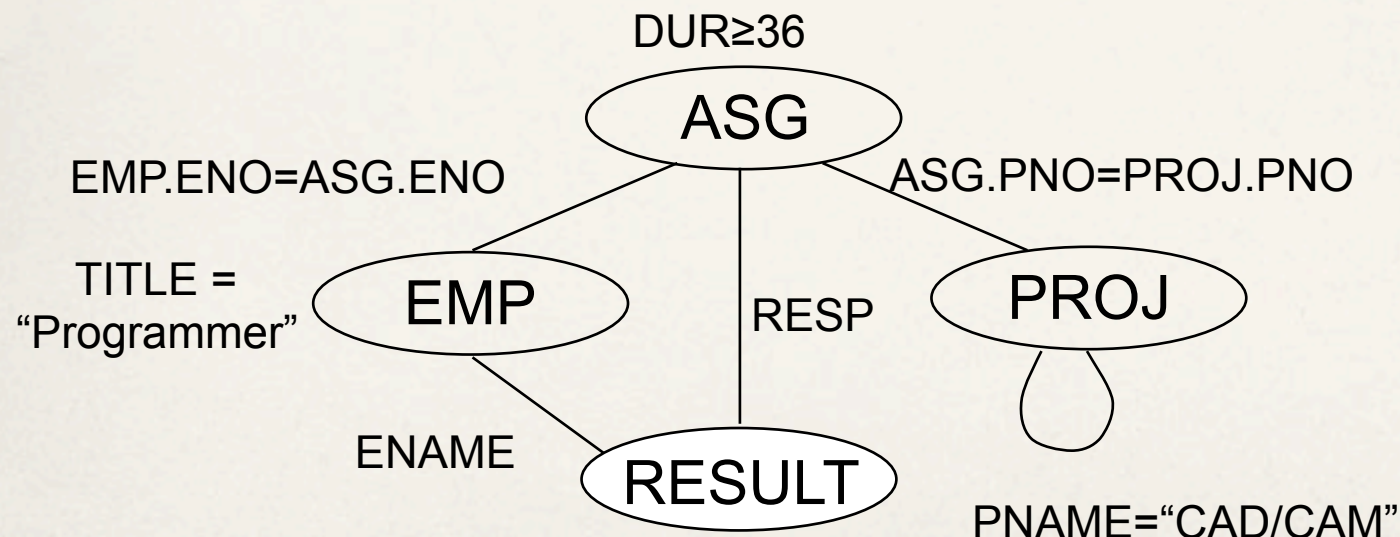
Analysis

- Refute incorrect queries
- Type incorrect
 - ➔ If any of its attribute or relation names are not defined in the global schema
 - ➔ If operations are applied to attributes of the wrong type
- Semantically incorrect
 - ➔ Components do not contribute in any way to the generation of the result
 - ➔ Only a subset of relational calculus queries can be tested for correctness
 - ➔ Those that do not contain disjunction and negation
 - ➔ To detect
 - ◆ connection graph (query graph)
 - ◆ join graph

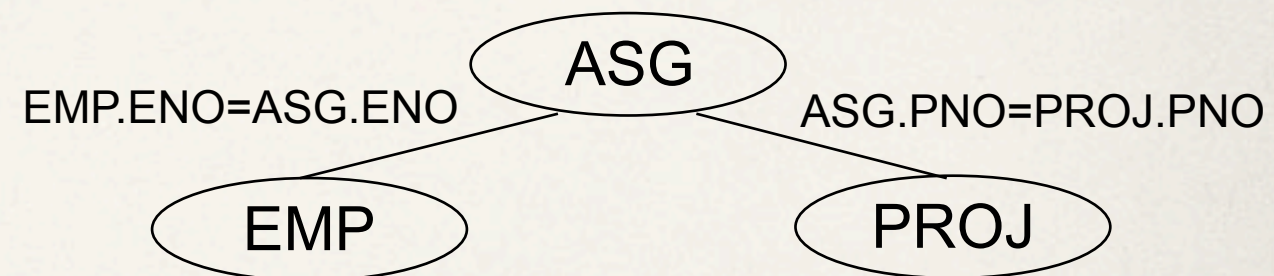
Analysis – Example

```
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"
```

Query graph



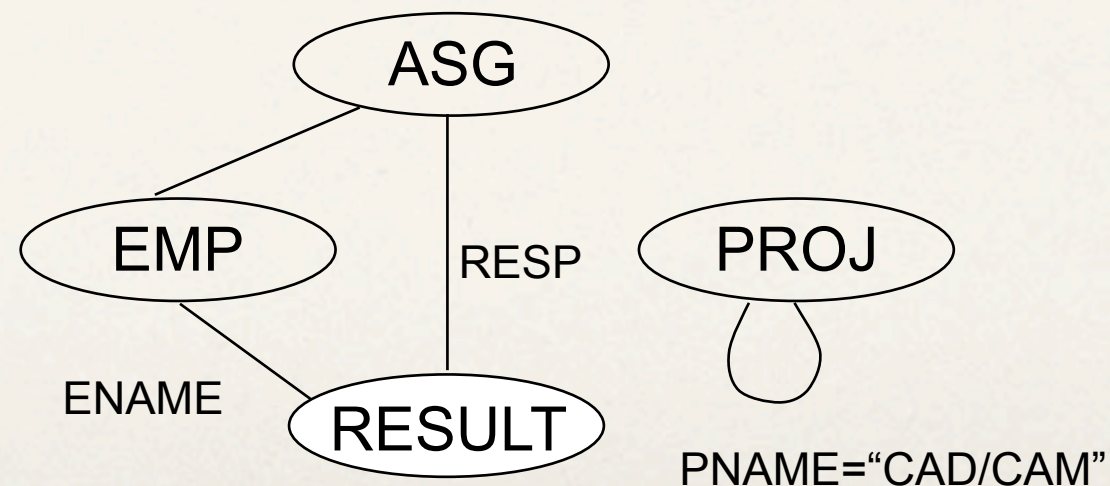
Join graph



Analysis

If the query graph is not connected, the query may be wrong or use Cartesian product

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



Simplification

- Why simplify?
 - ➔ Remember the example
- How? Use transformation rules

- ➔ Elimination of redundancy

- ♦ idempotency rules

$$p_1 \wedge \neg(p_1) \Leftrightarrow \text{false}$$

$$p_1 \wedge (p_1 \vee p_2) \Leftrightarrow p_1$$

$$p_1 \wedge \text{false} \Leftrightarrow p_1$$

...

- ➔ Application of transitivity
- ➔ Use of integrity rules

Simplification – Example

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



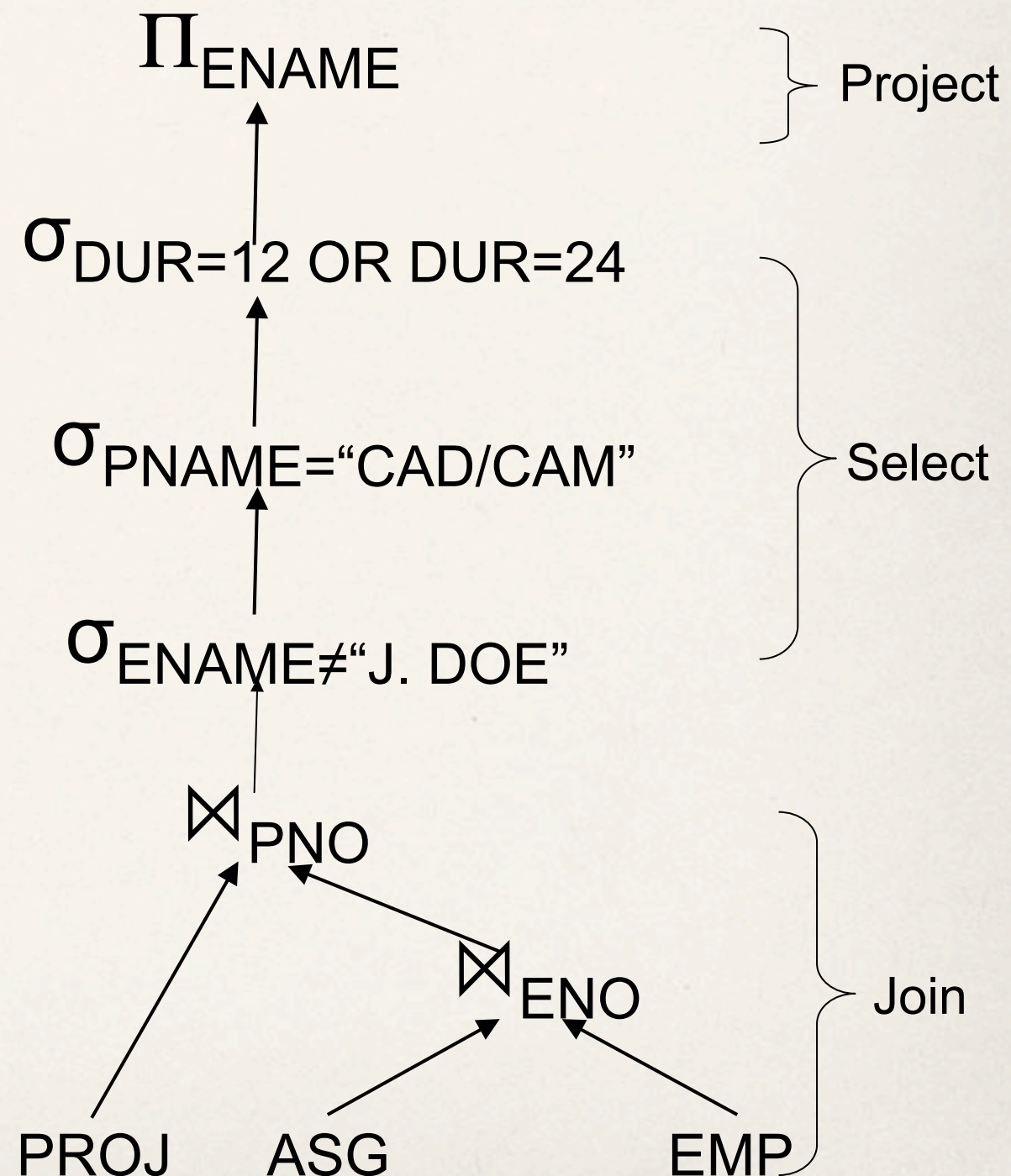
```
SELECT      TITLE
FROM        EMP
WHERE       EMP.ENAME = "J. Doe"
```


Restructuring

- Convert relational calculus to relational algebra
- Make use of query trees
- Example

Find the names of employees other than J. Doe who worked on the CAD/CAM project for either 1 or 2 years.

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



Restructuring – Transformation Rules

- Commutativity of binary operations
 - $R \times S \Leftrightarrow S \times R$
 - $R \bowtie S \Leftrightarrow S \bowtie R$
 - $R \cup S \Leftrightarrow S \cup R$
- Associativity of binary operations
 - $(R \times S) \times T \Leftrightarrow R \times (S \times T)$
 - $(R \bowtie S) \bowtie T \Leftrightarrow R \bowtie (S \bowtie T)$
- Idempotence of unary operations
 - $\Pi_{A'}(\Pi_{A'}(R)) \Leftrightarrow \Pi_{A'}(R)$
 - $\sigma_{p_1(A_1)}(\sigma_{p_2(A_2)}(R)) \Leftrightarrow \sigma_{p_1(A_1) \wedge p_2(A_2)}(R)$
where $R[A]$ and $A' \subseteq A, A'' \subseteq A$ and $A' \subseteq A''$
- Commuting selection with projection

Restructuring – Transformation Rules

- Commuting selection with binary operations

- $\sigma_{p(A)}(R \times S) \Leftrightarrow (\sigma_{p(A)}(R)) \times S$

- $\sigma_{p(A_i)}(R \bowtie_{(A_j B_k)} S) \Leftrightarrow (\sigma_{p(A_i)}(R)) \bowtie_{(A_j B_k)} S$

- $\sigma_{p(A_i)}(R \cup T) \Leftrightarrow \sigma_{p(A_i)}(R) \cup \sigma_{p(A_i)}(T)$

where A_i belongs to R and T

- Commuting projection with binary operations

- $\Pi_C(R \times S) \Leftrightarrow \Pi_{A'}(R) \times \Pi_{B'}(S)$

- $\Pi_C(R \bowtie_{(A_j B_k)} S) \Leftrightarrow \Pi_{A'}(R) \bowtie_{(A_j B_k)} \Pi_{B'}(S)$

- $\Pi_C(R \cup S) \Leftrightarrow \Pi_C(R) \cup \Pi_C(S)$

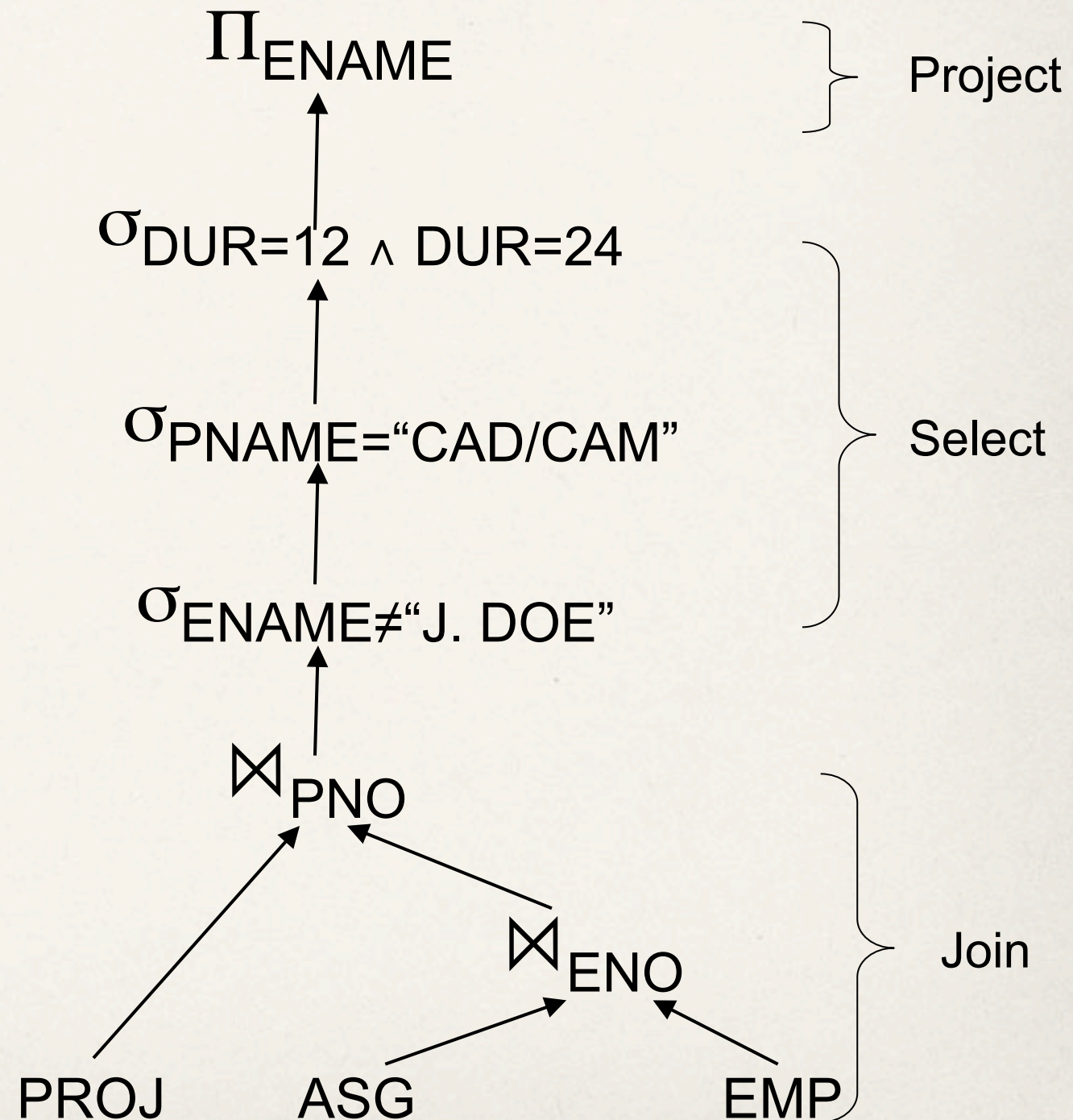
where $R[A]$ and $S[B]$; $C = A' \cup B'$ where $A' \subseteq A, B' \subseteq B$

Example

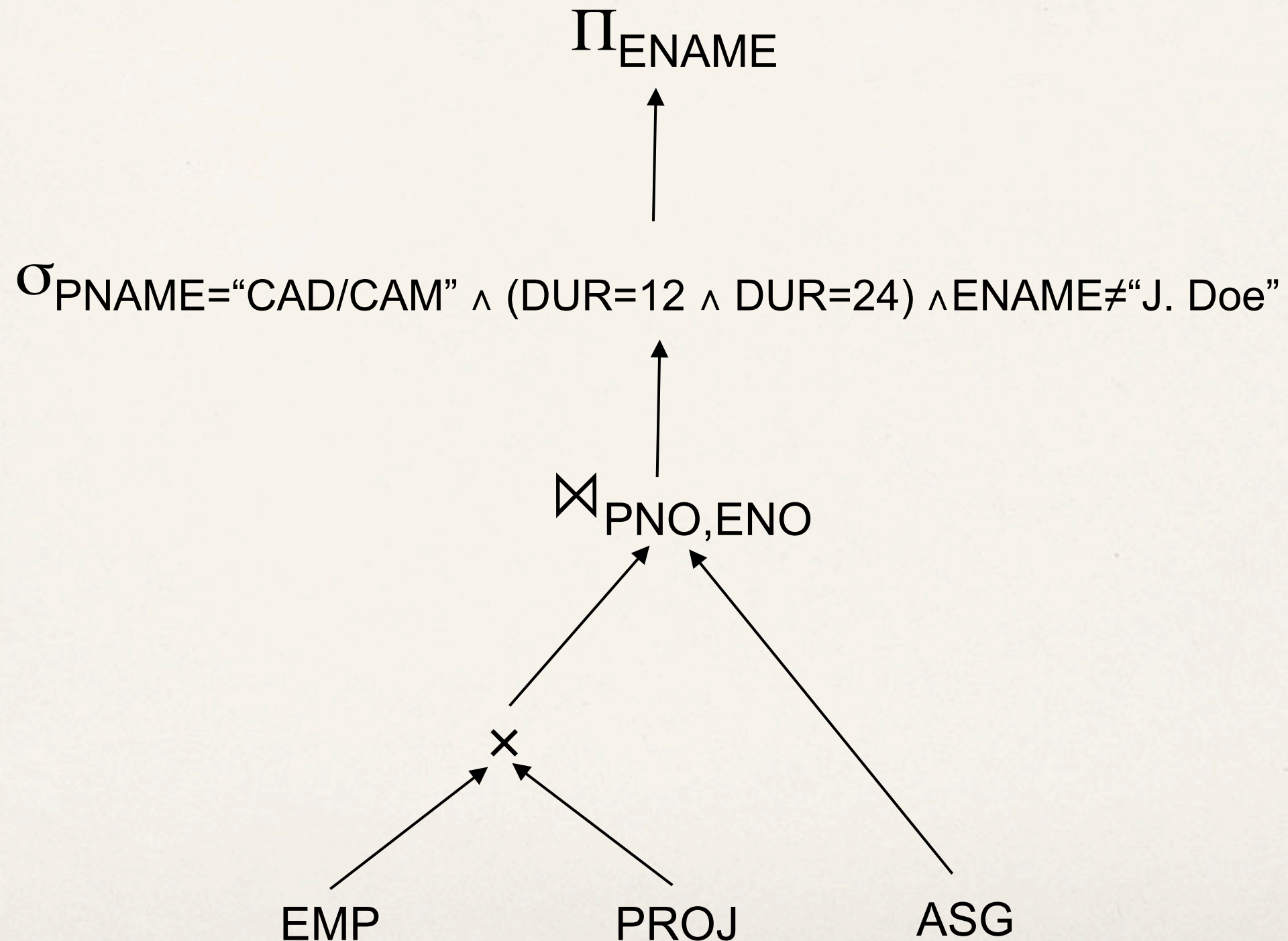
Recall the previous example:

Find the names of employees other than J. Doe who worked on the CAD/CAM project for either one or two years.

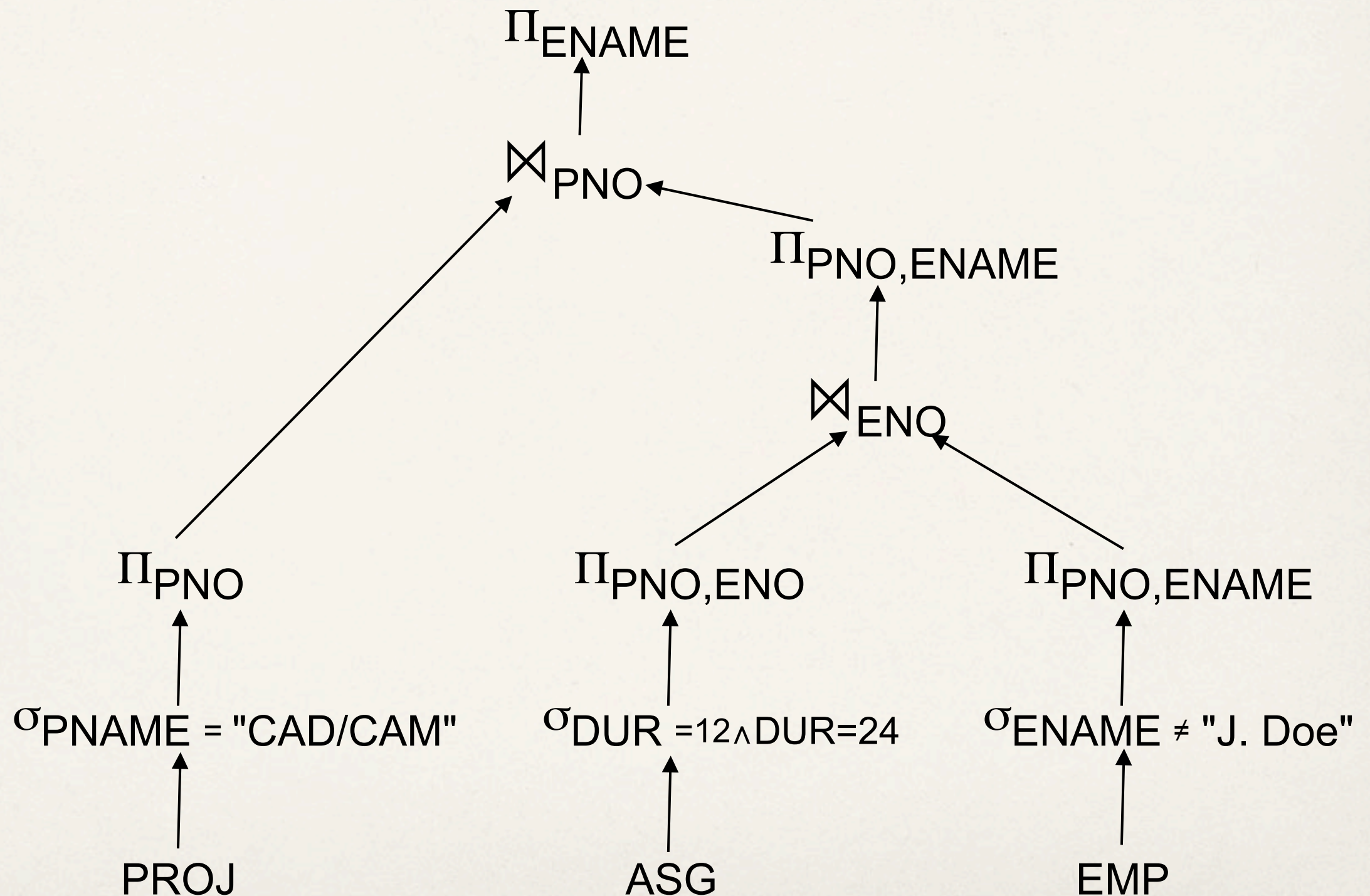
```
SELECT  ENAME
FROM    PROJ, ASG, EMP
WHERE   ASG.ENO=EMP.ENO
AND     ASG.PNO=PROJ.PNO
AND     ENAME ≠ "J. Doe"
AND     PROJ.PNAME="CAD/CAM"
AND     (DUR=12 OR DUR=24)
```



Equivalent Query



Restructuring



Step 2 – Data Localization

Input: Algebraic query on distributed relations

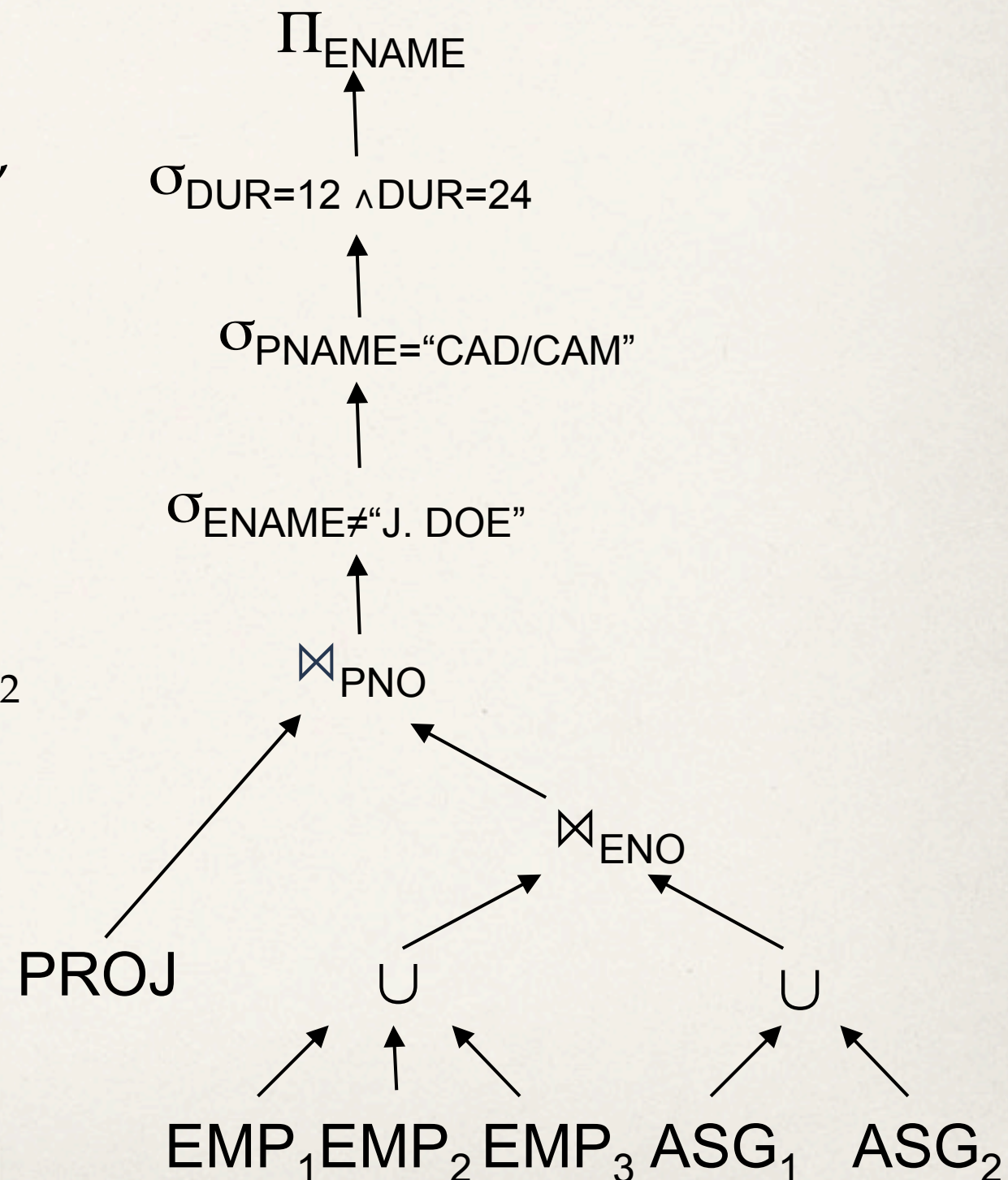
- Determine which fragments are involved
- **Localization program**
 - ➔ substitute for each global query its materialization program
 - ➔ optimize

Example

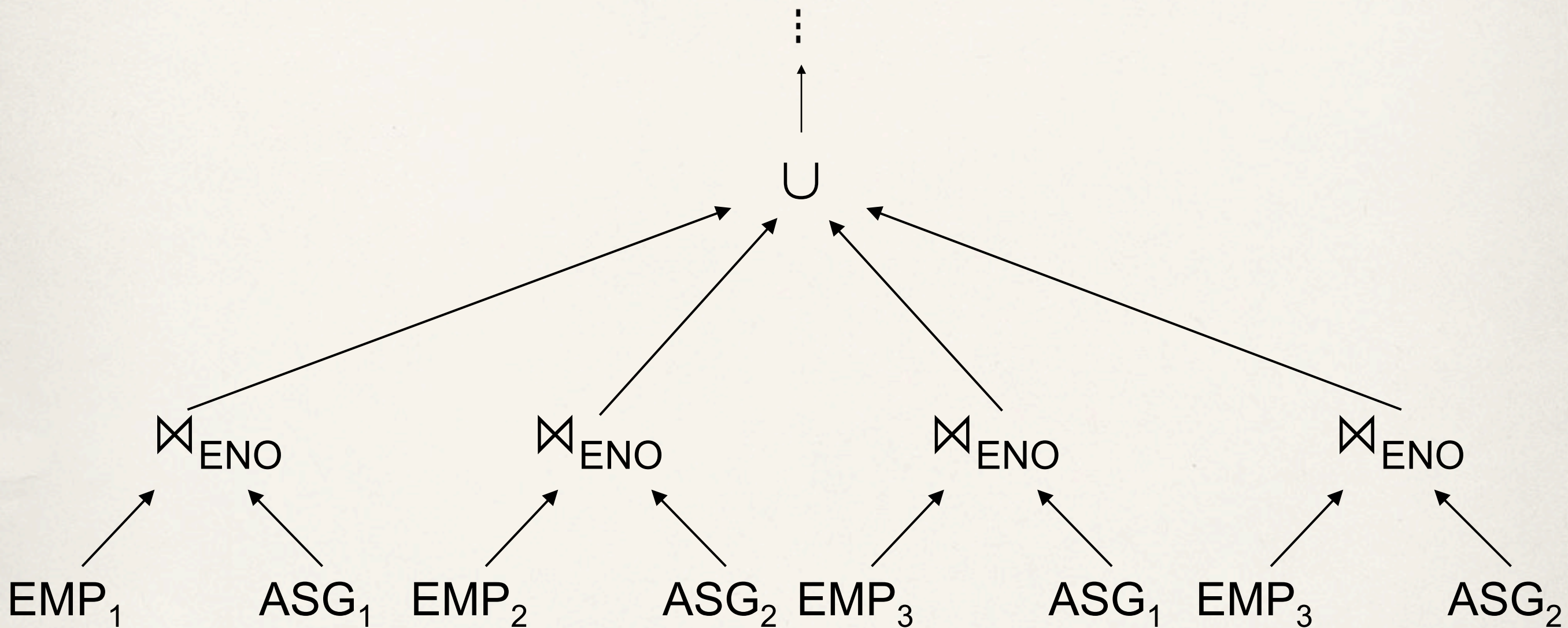
Assume

- EMP is fragmented into EMP_1 , EMP_2 , EMP_3 as follows:
 - ♦ $EMP_1 = \sigma_{ENO \leq "E3"}(EMP)$
 - ♦ $EMP_2 = \sigma_{"E3" < ENO \leq "E6"}(EMP)$
 - ♦ $EMP_3 = \sigma_{ENO \geq "E6"}(EMP)$
- ASG fragmented into ASG_1 and ASG_2 as follows:
 - ♦ $ASG_1 = \sigma_{ENO \leq "E3"}(ASG)$
 - ♦ $ASG_2 = \sigma_{ENO > "E3"}(ASG)$

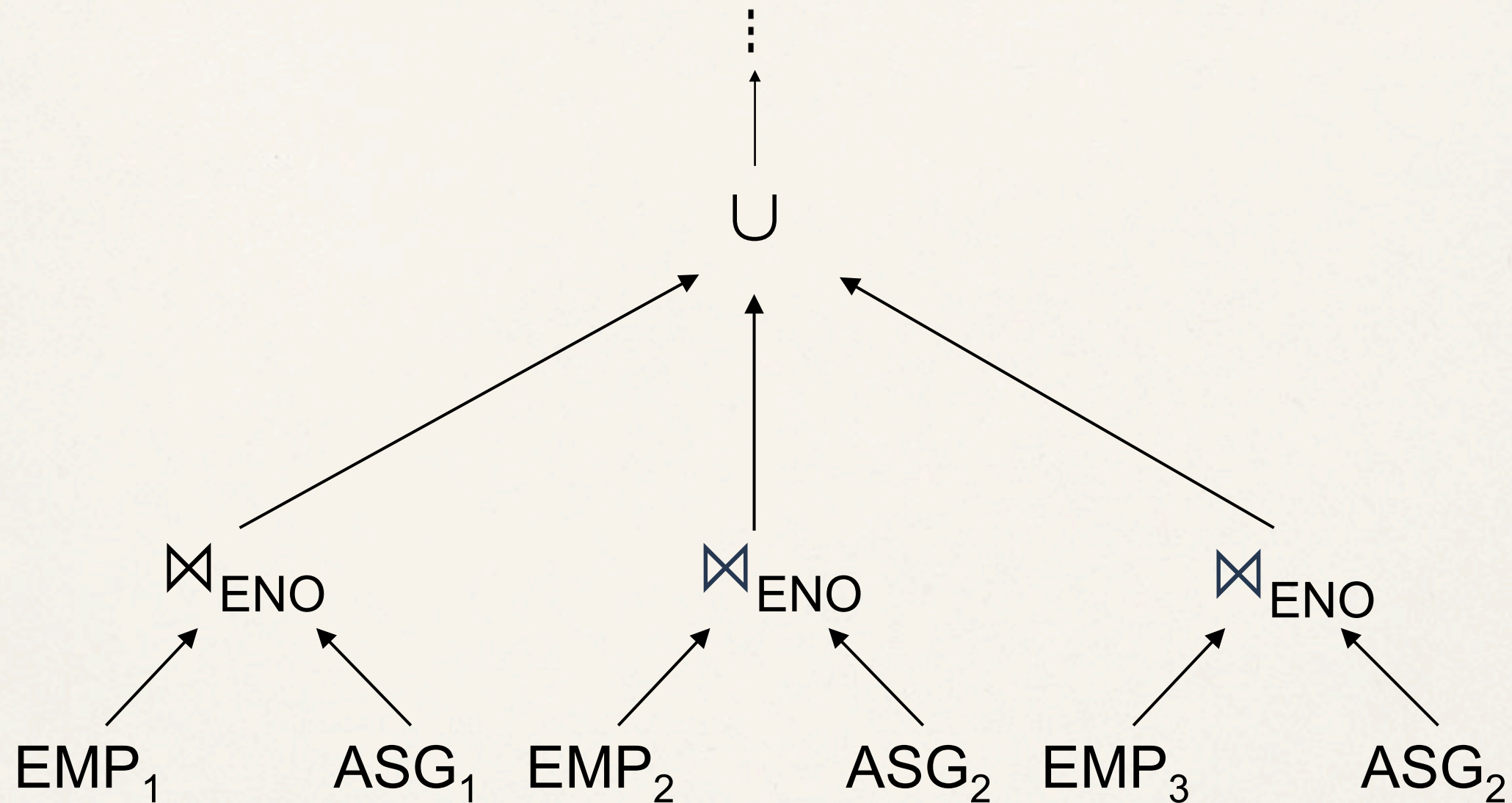
Replace EMP by $(EMP_1 \cup EMP_2 \cup EMP_3)$
and ASG by $(ASG_1 \cup ASG_2)$ in any query



Provides Parallelism



Eliminates Unnecessary Work



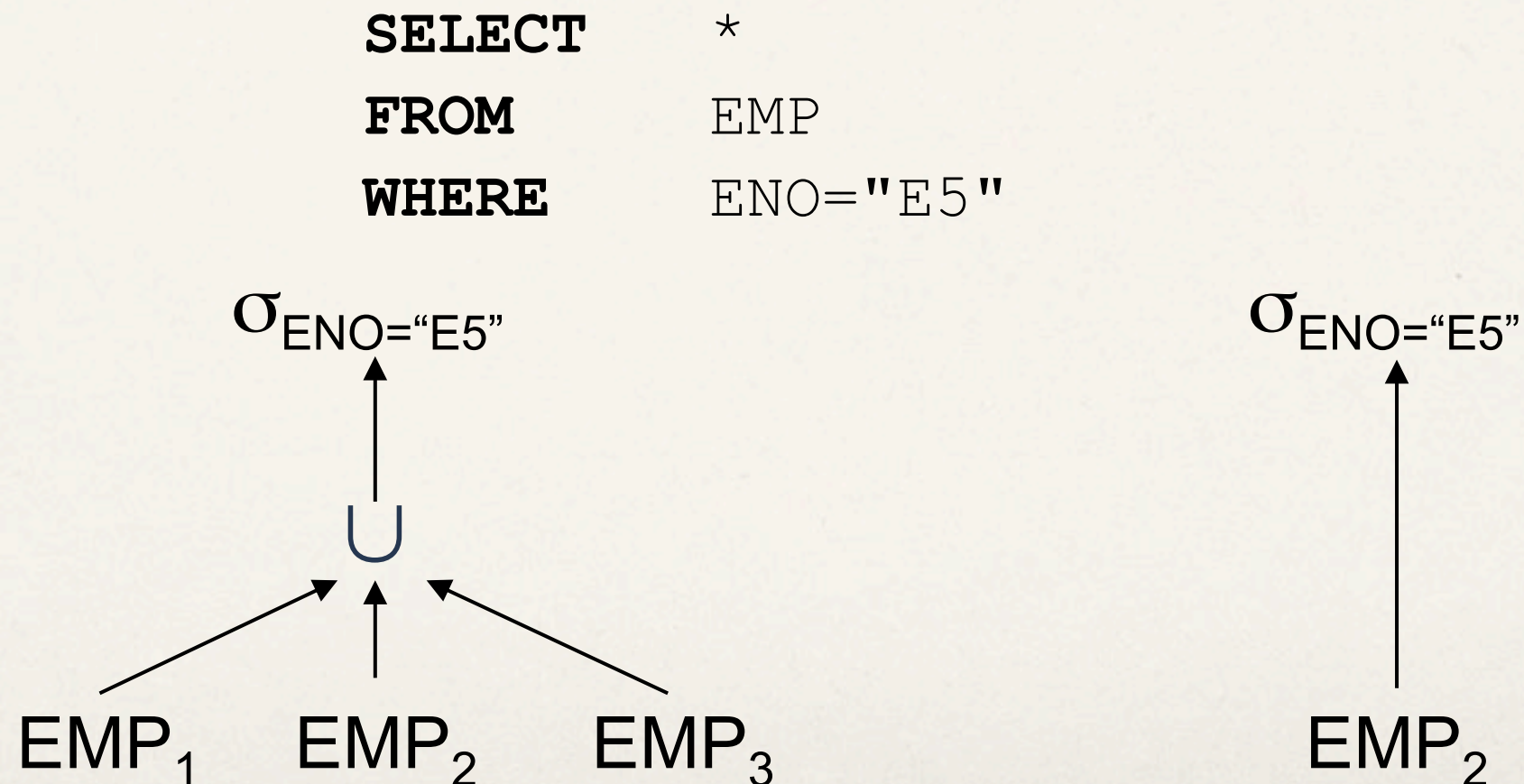
Reduction for PHF

- Reduction with selection

- Relation R and $F_R = \{R_1, R_2, \dots, R_w\}$ where $R_j = \sigma_{p_j}(R)$

$$\sigma_{p_i}(R_j) = \emptyset \text{ if } \forall x \text{ in } R: \neg(p_i(x) \wedge p_j(x))$$

- Example



Reduction for PHF

- Reduction with join

- ➔ Possible if fragmentation is done on join attribute
- ➔ Distribute join over union

$$(R_1 \cup R_2) \bowtie S \Leftrightarrow (R_1 \bowtie S) \cup (R_2 \bowtie S)$$

- ➔ Given $R_i = \sigma_{p_i}(R)$ and $R_j = \sigma_{p_j}(R)$

$$R_i \bowtie R_j = \emptyset \text{ if } \forall x \text{ in } R_i, \forall y \text{ in } R_j: \neg(p_i(x) \wedge p_j(y))$$

Reduction for PHF

- Assume EMP is fragmented as before and

→ $ASG_1: \sigma_{ENO \leq "E3"}(ASG)$

→ $ASG_2: \sigma_{ENO > "E3"}(ASG)$

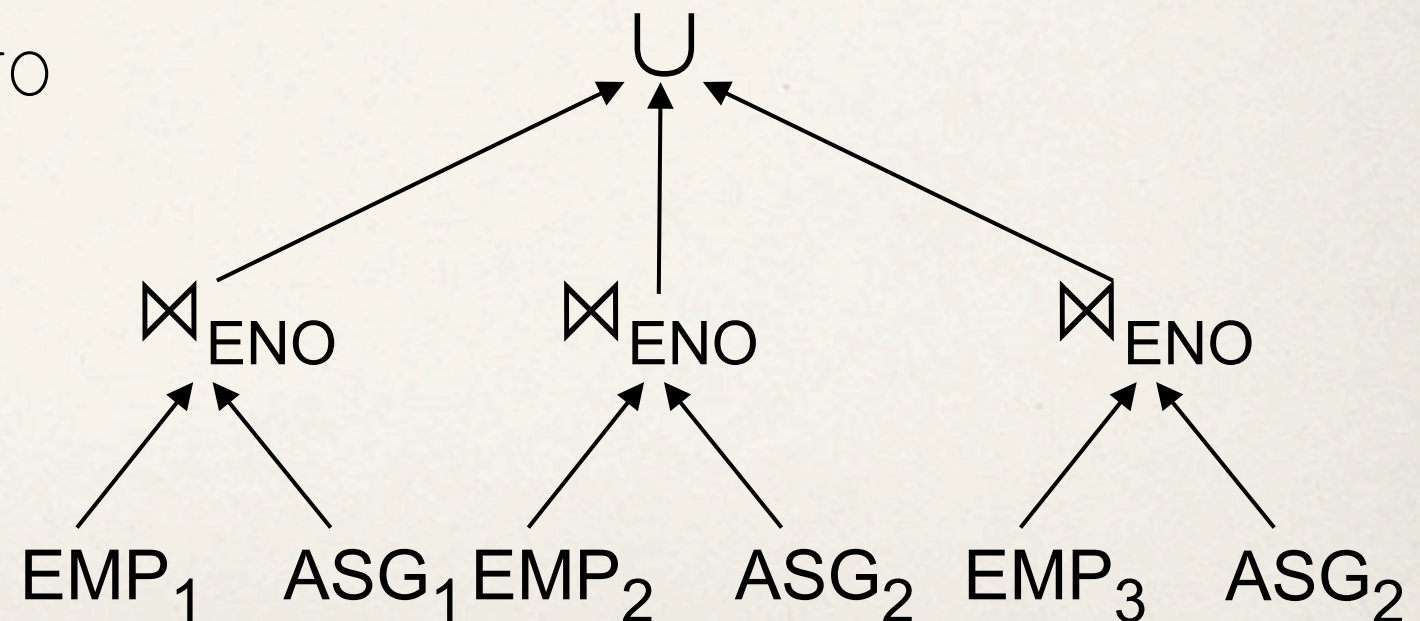
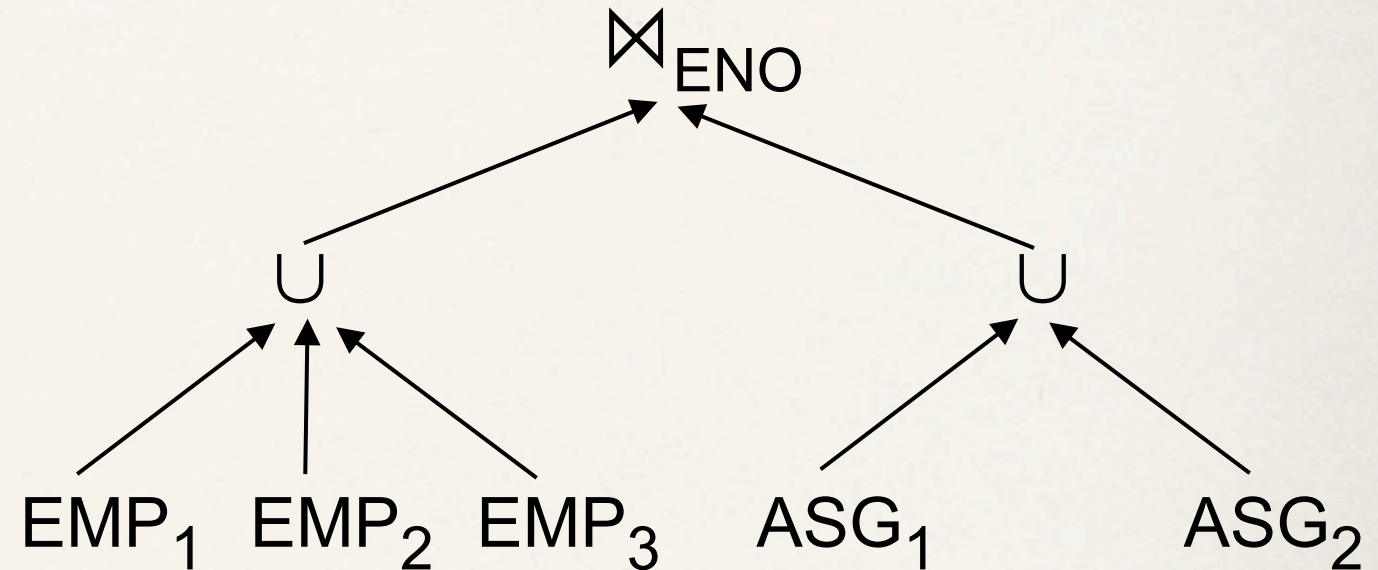
- Consider the query

```

SELECT      *
FROM        EMP, ASG
WHERE       EMP.ENO=ASG.ENO
    
```

- Distribute join over unions

- Apply the reduction rule



Reduction for VF

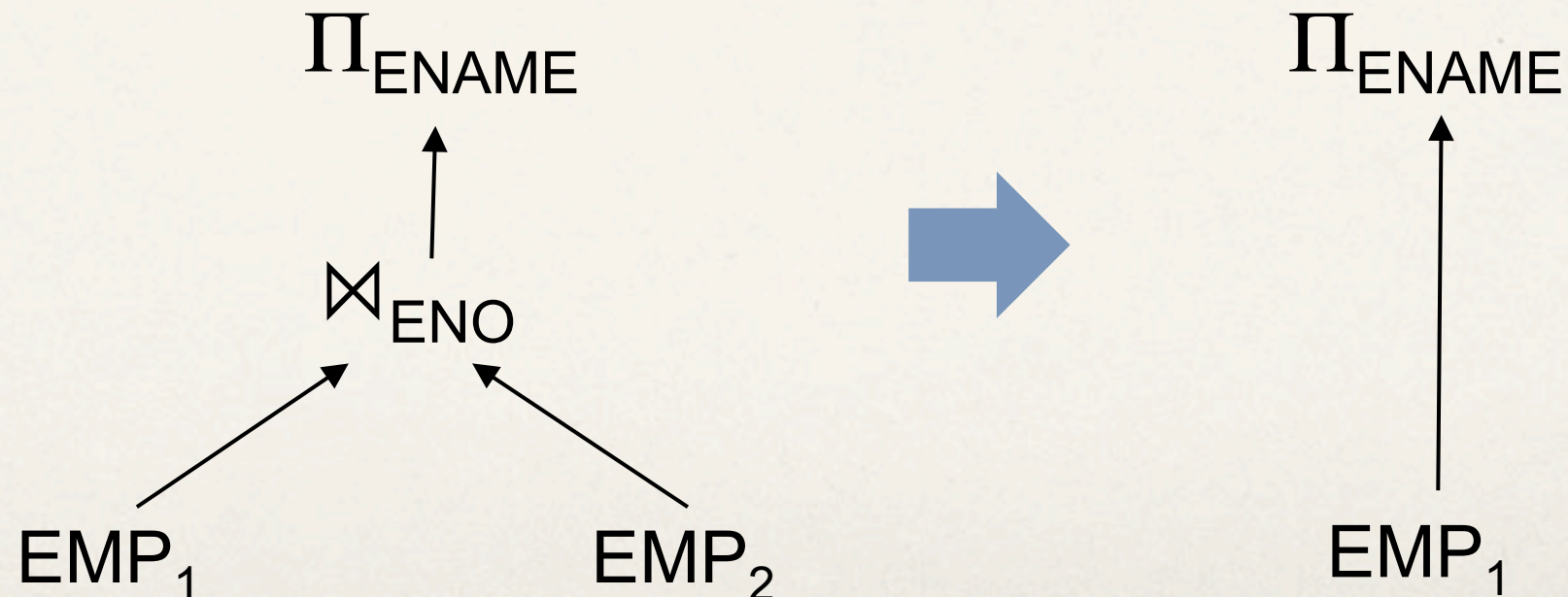
- Find useless (not empty) intermediate relations

Relation R defined over attributes $A = \{A_1, \dots, A_n\}$ vertically fragmented as $R_i = \Pi_{A'}(R)$ where $A' \subseteq A$:

$\Pi_{D,K}(R_i)$ is useless if the set of projection attributes D is not in A'

Example: $EMP_1 = \Pi_{ENO,ENAME}(EMP)$; $EMP_2 = \Pi_{ENO,TITLE}(EMP)$

SELECT ENAME
FROM EMP



Reduction for DHF

- Rule :

- Distribute joins over unions
- Apply the join reduction for horizontal fragmentation

- Example

$ASG_1: ASG \bowtie_{ENO} EMP_1$

$ASG_2: ASG \bowtie_{ENO} EMP_2$

$EMP_1: \sigma_{TITLE="Programmer"} (EMP)$

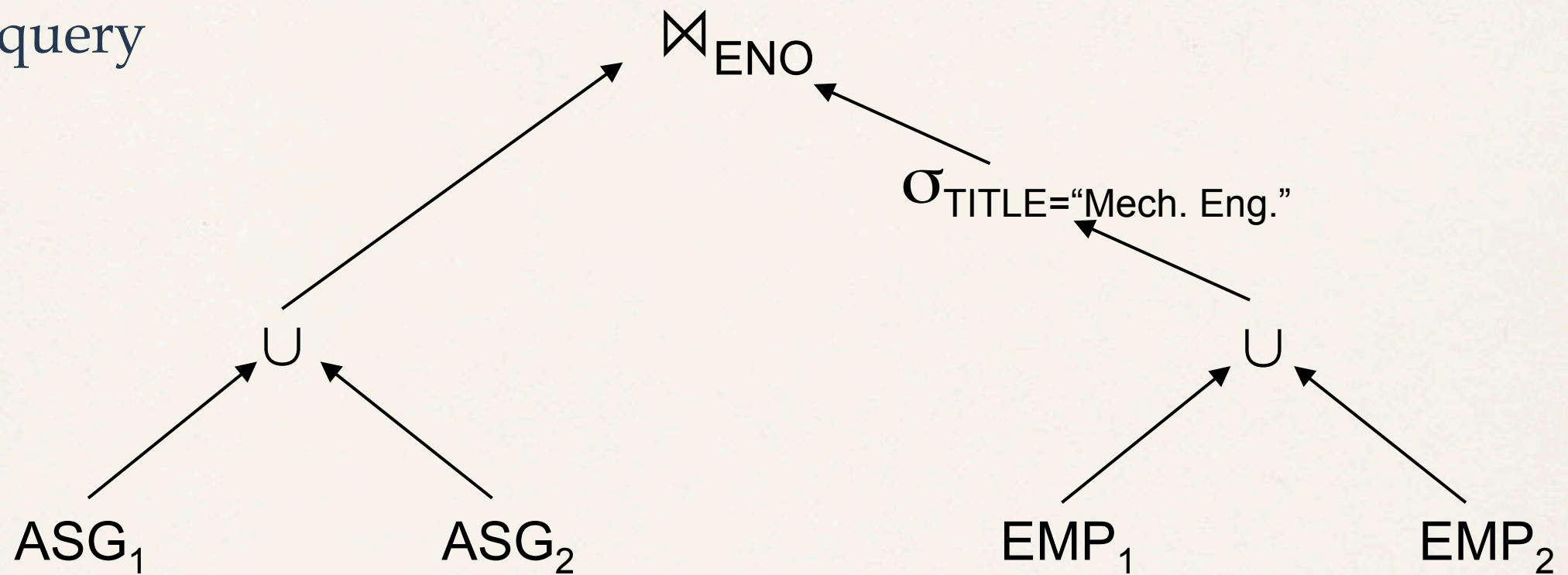
$EMP_2: \sigma_{TITLE="Programmer"} (EMP)$

- Query

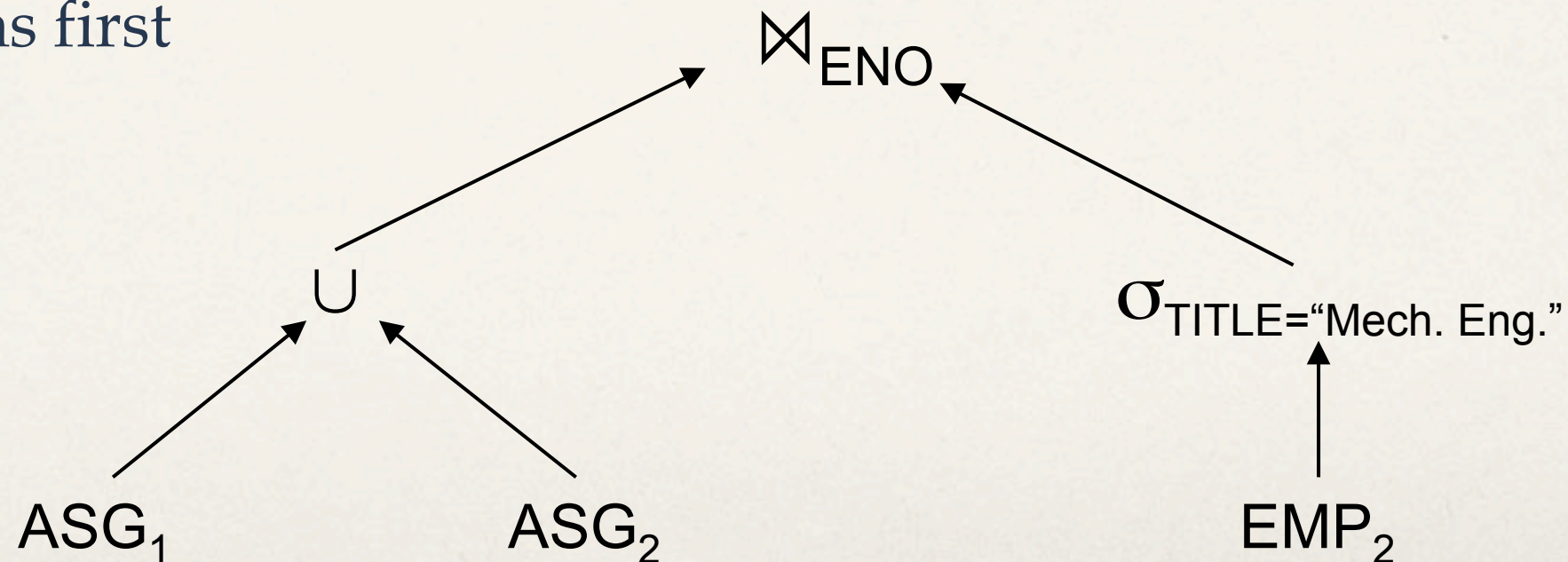
```
SELECT      *
FROM        EMP, ASG
WHERE       ASG.ENO = EMP.ENO
AND         EMP.TITLE = "Mech. Eng."
```


Reduction for DHF

Generic query

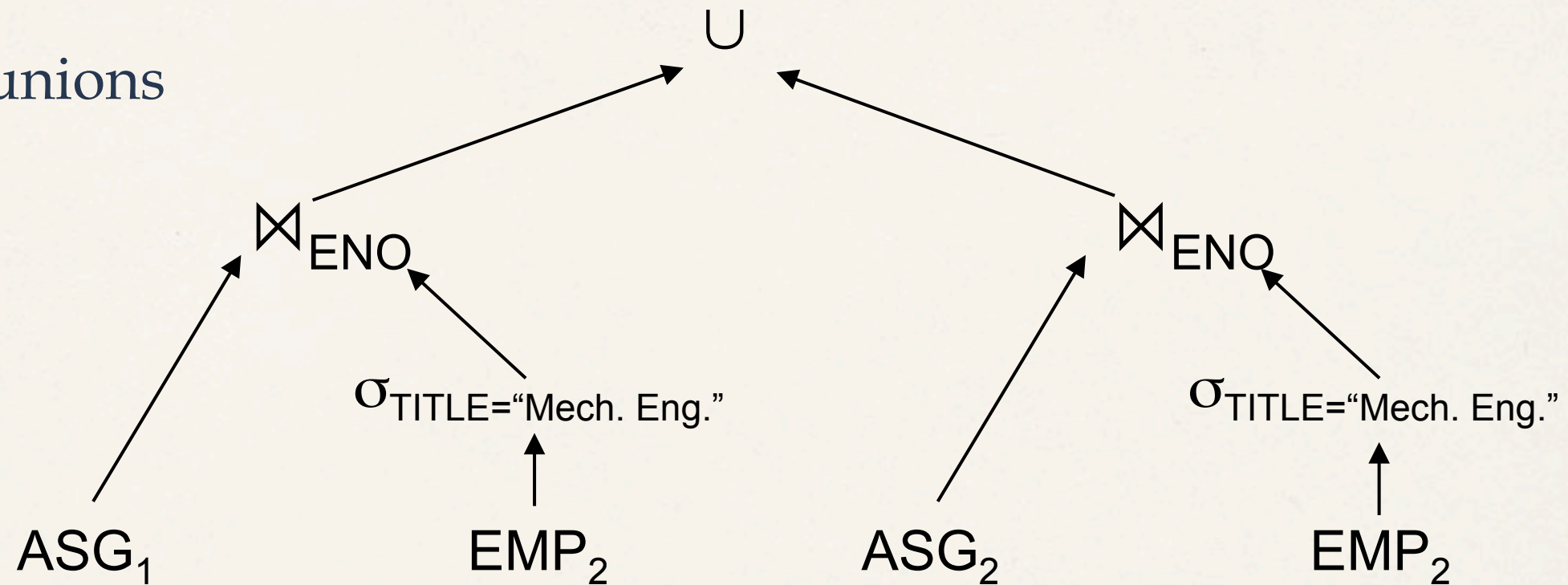


Selections first

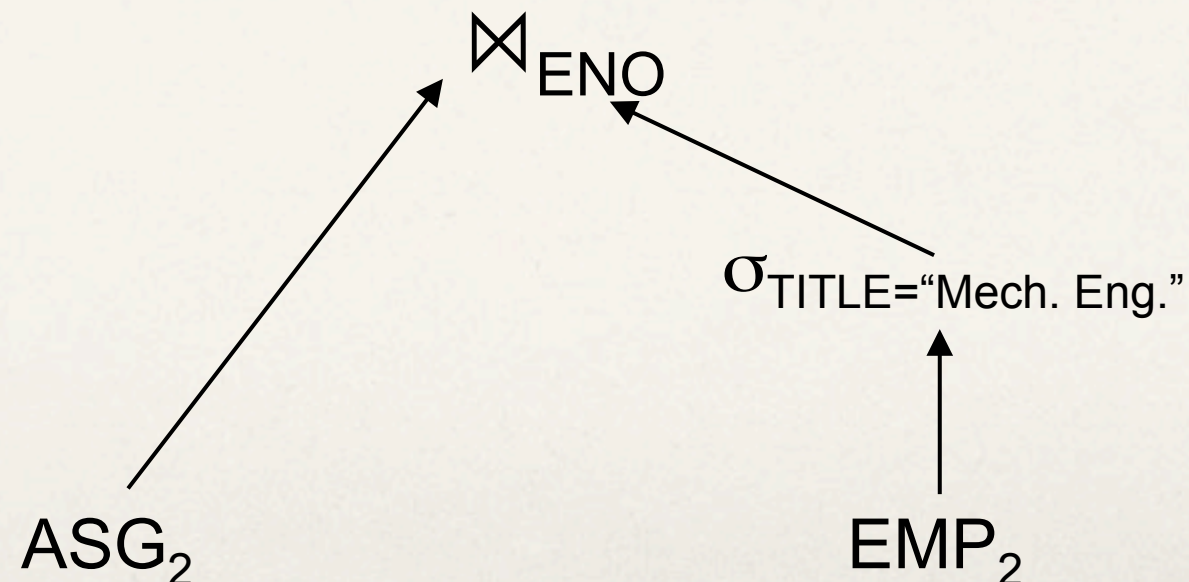


Reduction for DHF

Joins over unions



Elimination of the empty intermediate relations
(left sub-tree)



Reduction for Hybrid Fragmentation

- Combine the rules already specified:
 - ➔ Remove **empty relations** generated by contradicting selections on horizontal fragments;
 - ➔ Remove **useless relations** generated by projections on vertical fragments;
 - ➔ Distribute **joins over unions** in order to isolate and remove useless joins.

Reduction for HF

Example

Consider the following hybrid fragmentation:

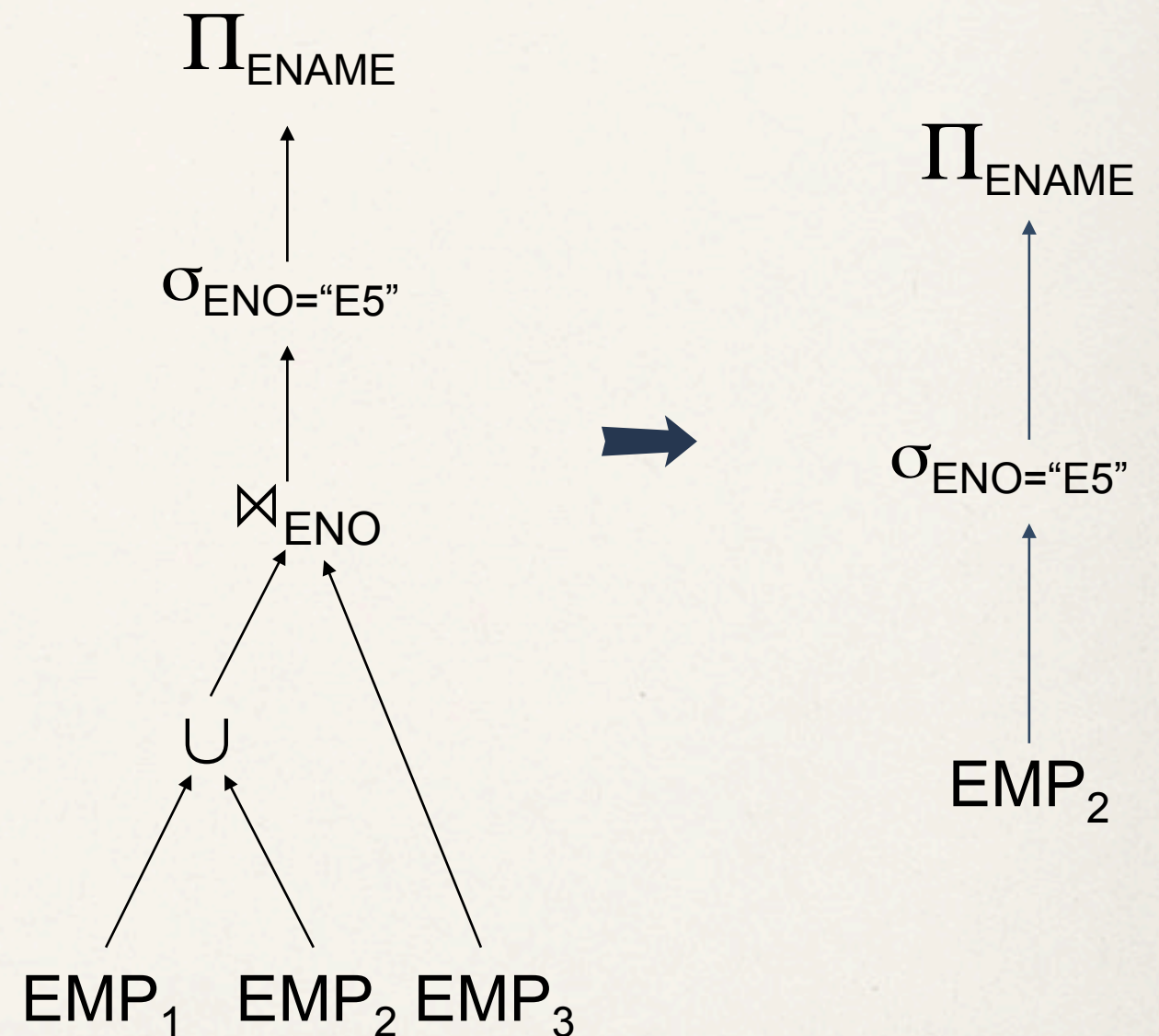
$$EMP_1 = \sigma_{ENO \leq "E4"} (\Pi_{ENO, ENAME} (EMP))$$

$$EMP_2 = \sigma_{ENO > "E4"} (\Pi_{ENO, ENAME} (EMP))$$

$$EMP_3 = \sigma_{ENO, TITLE} (EMP)$$

and the query

```
SELECT ENAME
FROM   EMP
WHERE  ENO = "E5"
```



Outline

- Introduction
- Background
- Distributed Database Design
- Database Integration
- Semantic Data Control
- Distributed Query Processing
 - ➔ Overview
 - ➔ Query decomposition and localization
 - ➔ Distributed query optimization
- Multidatabase Query Processing
- Distributed Transaction Management
- Data Replication
- Parallel Database Systems
- Distributed Object DBMS
- Peer-to-Peer Data Management
- Web Data Management
- Current Issues

Step 3 – Global Query Optimization

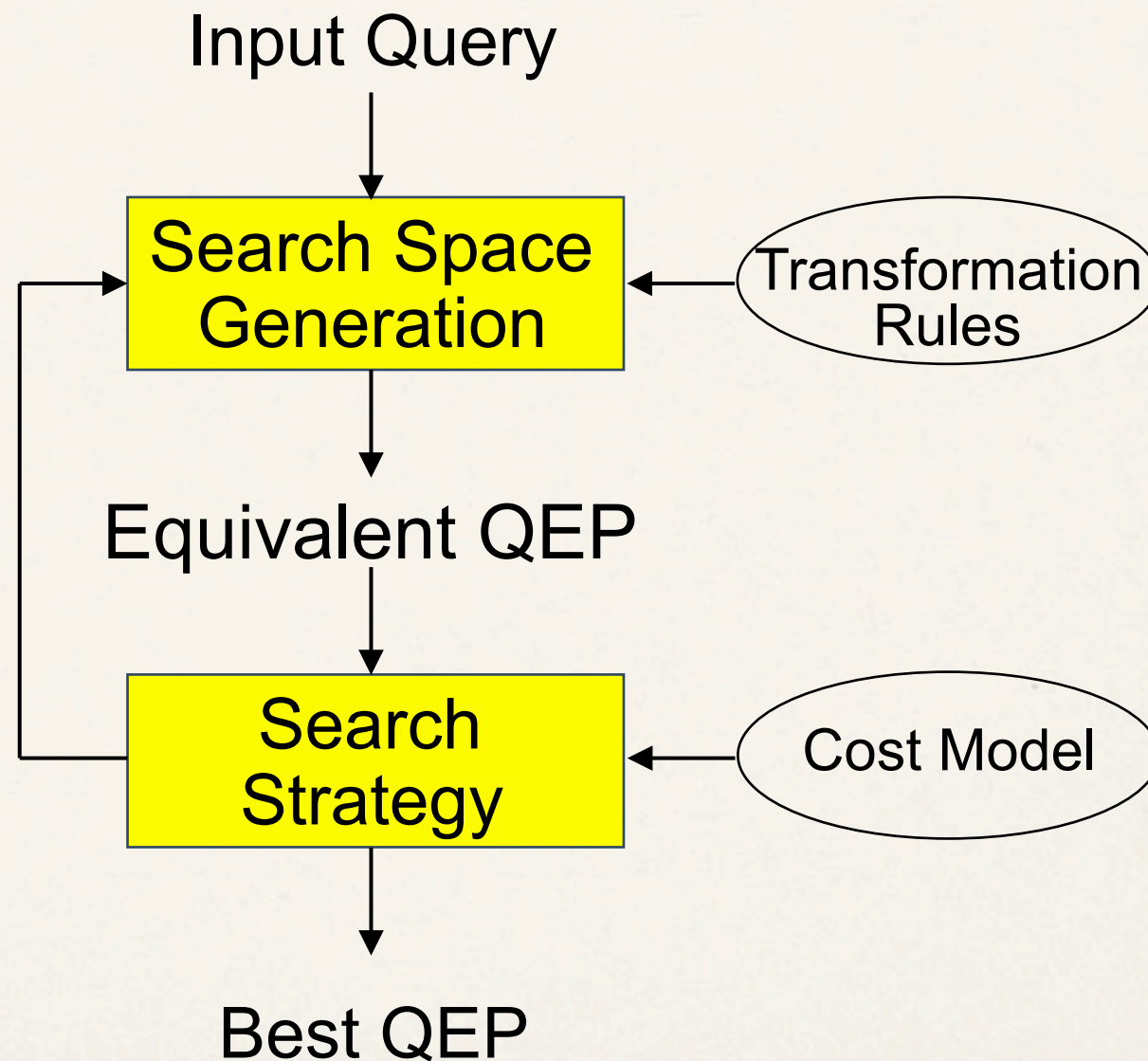
Input: Fragment query

- Find the *best* (not necessarily optimal) global schedule
 - ➔ Minimize a cost function
 - ➔ Distributed join processing
 - ◆ Bushy vs. linear trees
 - ◆ Which relation to ship where?
 - ◆ Ship-whole vs ship-as-needed
 - ➔ Decide on the use of semijoins
 - ◆ Semijoin saves on communication at the expense of more local processing.
 - ➔ Join methods
 - ◆ nested loop vs ordered joins (merge join or hash join)

Cost-Based Optimization

- Solution space
 - ➔ The set of equivalent algebra expressions (query trees).
- Cost function (in terms of time)
 - ➔ I/O cost + CPU cost + communication cost
 - ➔ These might have different weights in different distributed environments (LAN vs WAN).
 - ➔ Can also maximize throughput
- Search algorithm
 - ➔ How do we move inside the solution space?
 - ➔ Exhaustive search, heuristic algorithms (iterative improvement, simulated annealing, genetic,...)

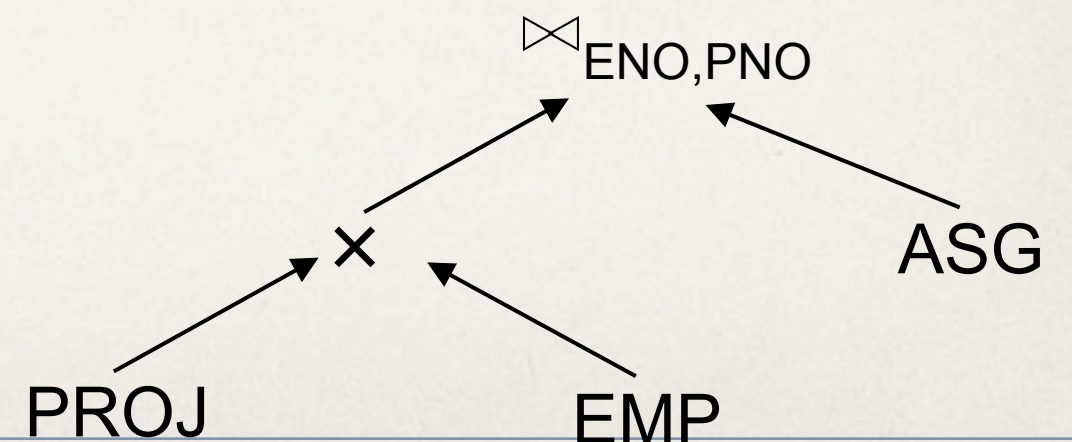
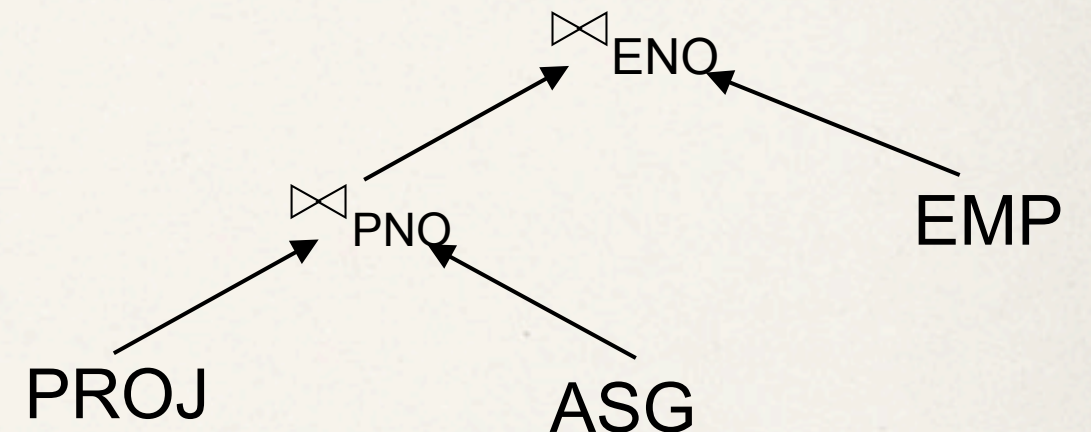
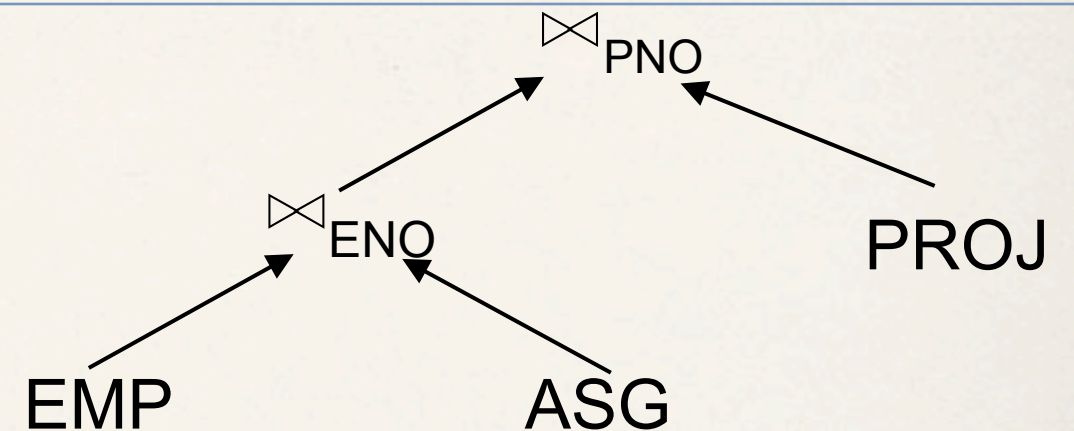
Query Optimization Process



Search Space

- Search space characterized by alternative execution
- Focus on join trees
- For N relations, there are $O(N!)$ equivalent join trees that can be obtained by applying commutativity and associativity rules

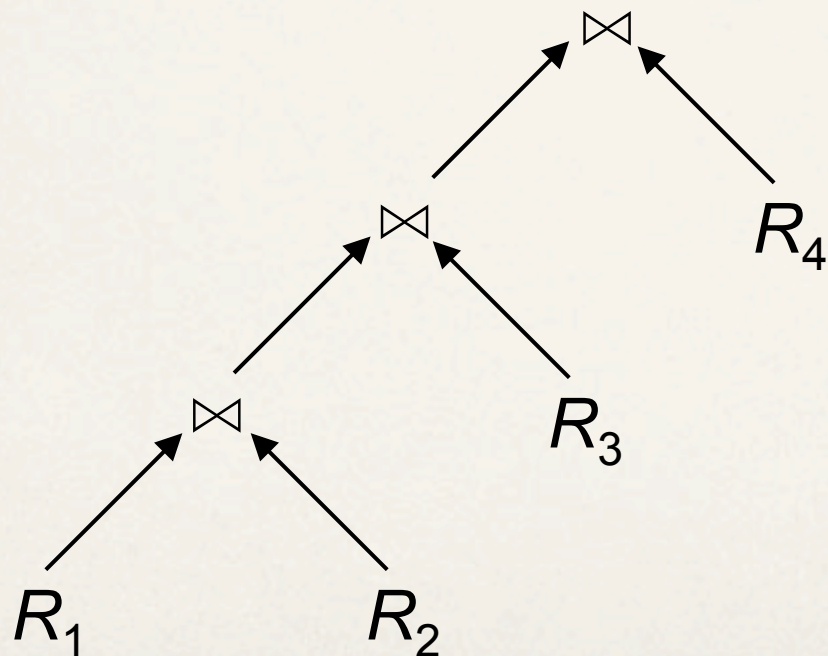
SELECT ENAME, RESP
FROM EMP, ASG, PROJ
WHERE EMP.ENO=ASG.ENO
AND ASG.PNO=PROJ.PNO



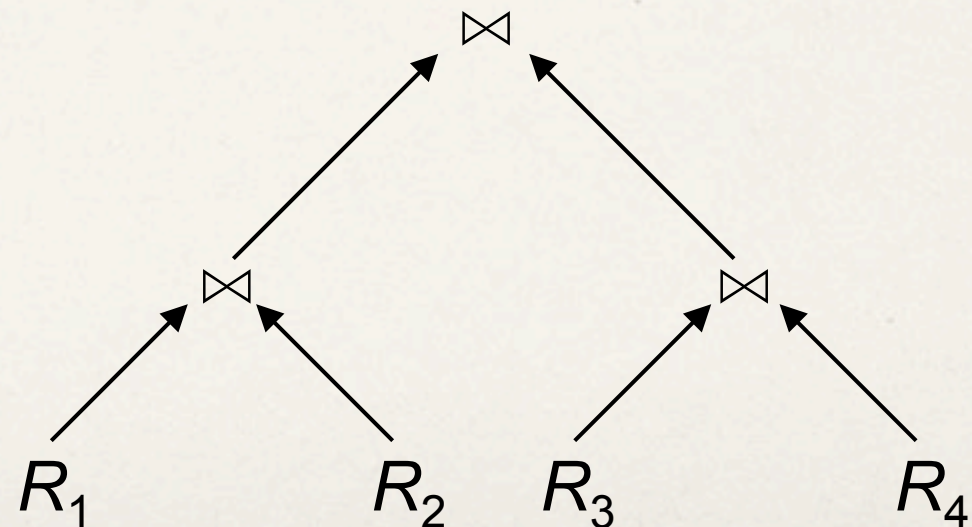
Search Space

- Restrict by means of heuristics
 - ➔ Perform unary operations before binary operations
 - ➔ ...
- Restrict the shape of the join tree
 - ➔ Consider only linear trees, ignore bushy ones

Linear Join Tree



Bushy Join Tree

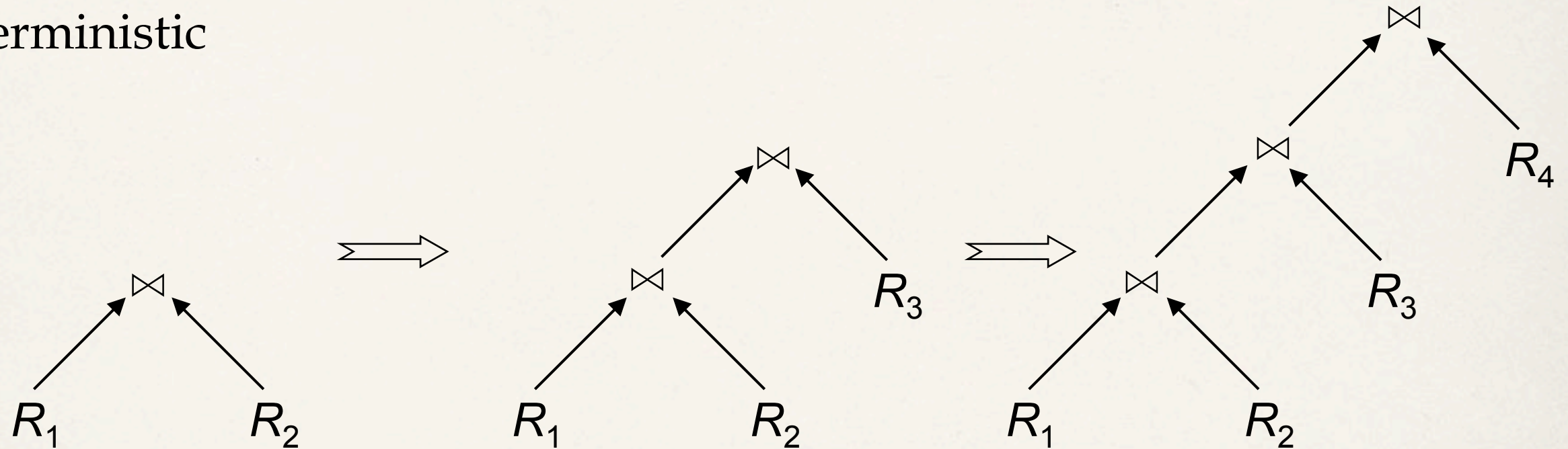


Search Strategy

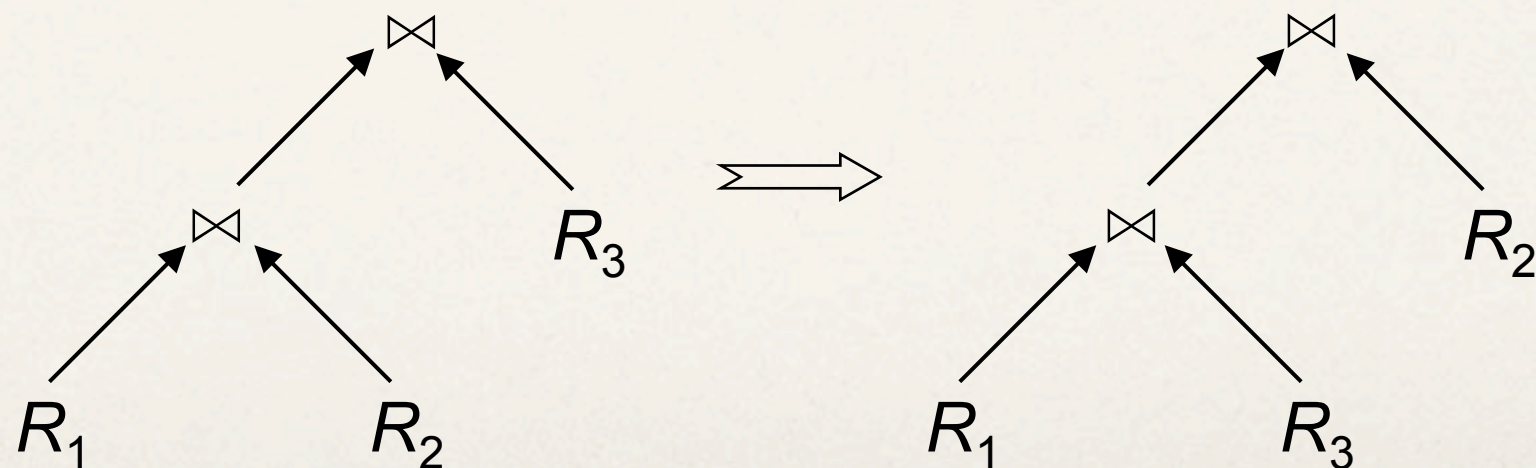
- How to “move” in the search space.
- Deterministic
 - ➔ Start from base relations and build plans by adding one relation at each step
 - ➔ Dynamic programming: breadth-first
 - ➔ Greedy: depth-first
- Randomized
 - ➔ Search for optimalities around a particular starting point
 - ➔ Trade optimization time for execution time
 - ➔ Better when > 10 relations
 - ➔ Simulated annealing
 - ➔ Iterative improvement

Search Strategies

- Deterministic



- Randomized



Cost Functions

- Total Time (or Total Cost)

- ➔ Reduce each cost (in terms of time) component individually
- ➔ Do as little of each cost component as possible
- ➔ Optimizes the utilization of the resources



Increases system throughput

- Response Time

- ➔ Do as many things as possible in parallel
- ➔ May increase total time because of increased total activity

Total Cost

Summation of all cost factors

Total cost = CPU cost + I/O cost + communication cost

CPU cost = unit instruction cost * no.of instructions

I/O cost = unit disk I/O cost * no. of disk I/Os

communication cost = message initiation + transmission

Total Cost Factors

- Wide area network
 - ➔ Message initiation and transmission costs high
 - ➔ Local processing cost is low (fast mainframes or minicomputers)
 - ➔ Ratio of communication to I/O costs = 20:1
- Local area networks
 - ➔ Communication and local processing costs are more or less equal
 - ➔ Ratio = 1:1.6

Response Time

Elapsed time between the initiation and the completion of a query

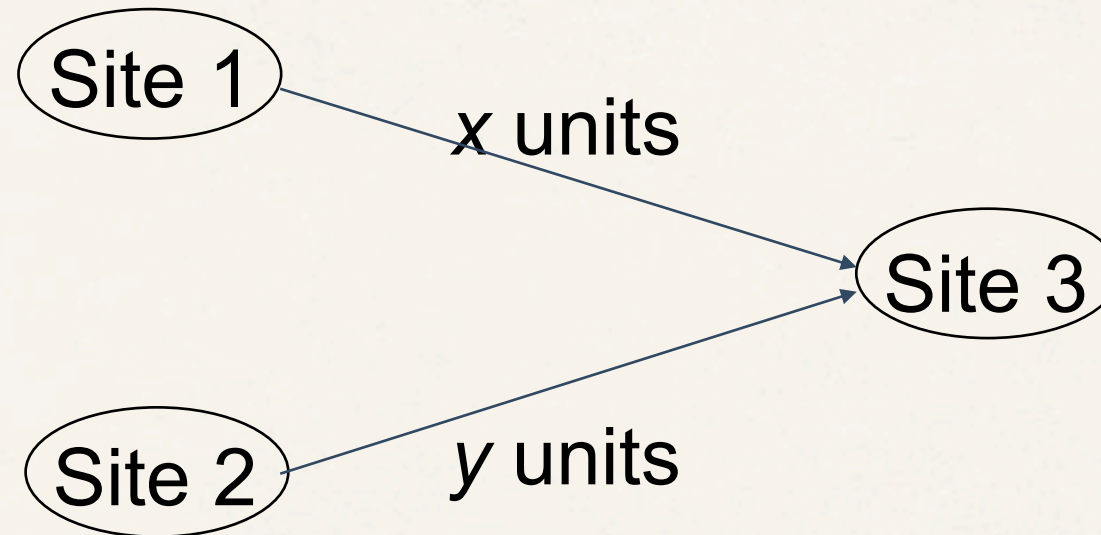
Response time = CPU time + I/O time + communication time

CPU time = unit instruction time * no. of sequential instructions

I/O time = unit I/O time * no. of sequential I/Os

communication time = unit msg initiation time * no. of sequential msg
+ unit transmission time * no. of sequential bytes

Example



Assume that only the communication cost is considered

Total time = $2 \cdot \text{message initialization time} + \text{unit transmission time} * (x+y)$

Response time = $\max \{ \text{time to send } x \text{ from 1 to 3, time to send } y \text{ from 2 to 3} \}$

time to send x from 1 to 3 = message initialization time
+ unit transmission time * x

time to send y from 2 to 3 = message initialization time
+ unit transmission time * y

Optimization Statistics

- Primary cost factor: **size of intermediate relations**
 - ➔ Need to estimate their sizes
- Make them precise \Rightarrow more costly to maintain
- Simplifying assumption: uniform distribution of attribute values in a relation

Statistics

- For each relation $R[A_1, A_2, \dots, A_n]$ fragmented as R_1, \dots, R_r
 - ➔ length of each attribute: $length(A_i)$
 - ➔ the number of distinct values for each attribute in each fragment: $card(\Pi_{A_i} R_j)$
 - ➔ maximum and minimum values in the domain of each attribute: $min(A_i), max(A_i)$
 - ➔ the cardinalities of each domain: $card(dom[A_i])$
- The cardinalities of each fragment: $card(R_j)$ Selectivity factor of each operation for relations
 - ➔ For joins

$$SF_{\bowtie}(R, S) = \frac{card(R \bowtie S)}{card(R) * card(S)}$$

Intermediate Relation Sizes

Selection

$$size(R) = card(R) \cdot length(R)$$

$$card(\sigma_F(R)) = SF_\sigma(F) \cdot card(R)$$

where

$$SF_\sigma(A = value) = \frac{1}{card(\Pi_A(R))}$$

$$SF_\sigma(A > value) = \frac{max(A) - value}{max(A) - min(A)}$$

$$SF_\sigma(A < value) = \frac{value - max(A)}{max(A) - min(A)}$$

$$SF_\sigma(p(A_i) \wedge p(A_j)) = SF_\sigma(p(A_i)) \cdot SF_\sigma(p(A_j))$$

$$SF_\sigma(p(A_i) \vee p(A_j)) = SF_\sigma(p(A_i)) + SF_\sigma(p(A_j)) - (SF_\sigma(p(A_i)) \cdot SF_\sigma(p(A_j)))$$

$$SF_\sigma(A \in \{value\}) = SF_\sigma(A = value) * card(\{values\})$$

Intermediate Relation Sizes

Projection

$$\text{card}(\Pi_A(R)) = \text{card}(R)$$

Cartesian Product

$$\text{card}(R \cdot S) = \text{card}(R) * \text{card}(S)$$

Union

$$\text{upper bound: } \text{card}(R \cup S) = \text{card}(R) + \text{card}(S)$$

$$\text{lower bound: } \text{card}(R \cup S) = \max\{\text{card}(R), \text{card}(S)\}$$

Set Difference

$$\text{upper bound: } \text{card}(R - S) = \text{card}(R)$$

$$\text{lower bound: } 0$$

Intermediate Relation Size

Join

- ➔ Special case: A is a key of R and B is a foreign key of S

$$\text{card}(R \bowtie_{A=B} S) = \text{card}(S)$$

- ➔ More general:

$$\text{card}(R \bowtie S) = SF_{\bowtie} * \text{card}(R) \cdot \text{card}(S)$$

Semijoin

$$\text{card}(R \ltimes_A S) = SF_{\ltimes}(S.A) * \text{card}(R)$$

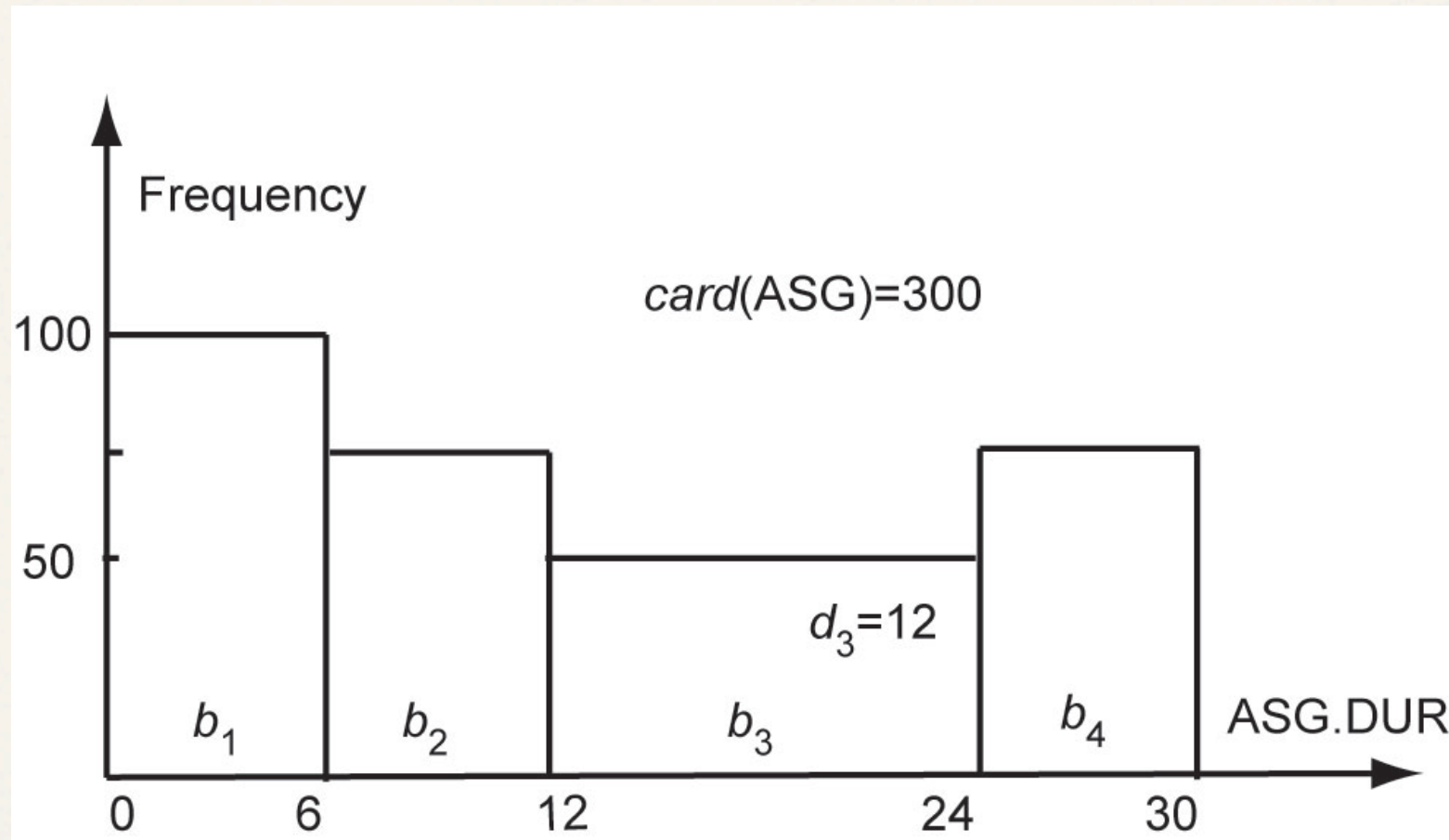
where

$$SF_{\ltimes}(R \ltimes_A S) = SF_{\ltimes}(S.A) = \frac{\text{card}(\Pi_A(S))}{\text{card}(\text{dom}[A])}$$

Histograms for Selectivity Estimation

- For skewed data, the uniform distribution assumption of attribute values yields inaccurate estimations
- Use an histogram for each skewed attribute A
 - ➔ Histogram = set of buckets
 - ◆ Each bucket describes a range of values of A , with its average frequency f (number of tuples with A in that range) and number of distinct values d
 - ◆ Buckets can be adjusted to different ranges
- Examples
 - ➔ Equality predicate
 - ◆ With (value in Range_i), we have: $SF_{\sigma}(A = \text{value}) = 1/d_i$
 - ➔ Range predicate
 - ◆ Requires identifying relevant buckets and summing up their frequencies

Histogram Example



For $ASG.DUR=18$: we have $SF=1/12$ so the card of selection is $300/12 = 25$ tuples

For $ASG.DUR \leq 18$: we have $\min(range_3)=12$ and $\max(range_3)=24$ so the card. of selection is $100+75+(((18-12)/(24-12))*50) = 200$ tuples

Centralized Query Optimization

- Dynamic (Ingres project at UCB)
 - ➔ Interpretive
- Static (System R project at IBM)
 - ➔ Exhaustive search
- Hybrid (Volcano project at OGI)
 - ➔ Choose node within plan

Dynamic Algorithm

- ① Decompose each multi-variable query into a sequence of mono-variable queries with a common variable
- ② Process each by a one variable query processor
 - ➔ Choose an initial execution plan (heuristics)
 - ➔ Order the rest by considering intermediate relation sizes



No statistical information is maintained

Dynamic Algorithm— Decomposition

- Replace an n variable query q by a series of queries

$$q_1 \rightarrow q_2 \rightarrow \dots \rightarrow q_n$$

where q_i uses the result of q_{i-1} .

- Detachment

→ Query q decomposed into $q' \rightarrow q''$ where q' and q'' have a common variable which is the result of q'

- Tuple substitution

→ Replace the value of each tuple with actual values and simplify the query

$$q(V_1, V_2, \dots, V_n) \rightarrow (q'(t_1, V_2, V_2, \dots, V_n), t_1 \in R)$$

Detachment

q : **SELECT** $V_2 \cdot A_2, V_3 \cdot A_3, \dots, V_n \cdot A_n$
 FROM $R_1 \ V_1, \dots, R_n \ V_n$
 WHERE $P_1 (V_1 \cdot A_1')$ **AND** $P_2 (V_1 \cdot A_1, V_2 \cdot A_2, \dots, V_n \cdot A_n)$



q' : **SELECT** $V_1 \cdot A_1$ **INTO** R_1'
 FROM $R_1 \ V_1$
 WHERE $P_1 (V_1 \cdot A_1)$

q'' : **SELECT** $V_2 \cdot A_2, \dots, V_n \cdot A_n$
 FROM $R_1' \ V_1, R_2 \ V_2, \dots, R_n \ V_n$
 WHERE $P_2 (V_1 \cdot A_1, V_2 \cdot A_2, \dots, V_n \cdot A_n)$

Detachment Example

Names of employees working on CAD/CAM project

q_1 : **SELECT** EMP.ENAME
FROM EMP, ASG, PROJ
WHERE EMP.ENO=ASG.ENO
AND ASG.PNO=PROJ.PNO
AND PROJ.PNAME="CAD/CAM"



q_{11} : **SELECT** PROJ.PNO **INTO** JVAR
FROM PROJ
WHERE PROJ.PNAME="CAD/CAM"

q' : **SELECT** EMP.ENAME
FROM EMP, ASG, JVAR
WHERE EMP.ENO=ASG.ENO
AND ASG.PNO=JVAR.PNO

Detachment Example (cont'd)

q' : **SELECT** EMP.ENAME
 FROM EMP, ASG, JVAR
 WHERE EMP.ENO=ASG.ENO
 AND ASG.PNO=JVAR.PNO



q_{12} : **SELECT** ASG.ENO **INTO** GVAR
 FROM ASG, JVAR
 WHERE ASG.PNO=JVAR.PNO

q_{13} : **SELECT** EMP.ENAME
 FROM EMP, GVAR
 WHERE EMP.ENO=GVAR.ENO

Tuple Substitution

q_{11} is a mono-variable query

q_{12} and q_{13} is subject to tuple substitution

Assume GVAR has two tuples only: $\langle E1 \rangle$ and $\langle E2 \rangle$

Then q_{13} becomes

q_{131} : **SELECT** EMP.ENAME
 FROM EMP
 WHERE EMP.ENO="E1"

q_{132} : **SELECT** EMP.ENAME
 FROM EMP
 WHERE EMP.ENO="E2"

Static Algorithm

- ① Simple (i.e., mono-relation) queries are executed according to the best access path
- ② Execute joins
 - ➔ Determine the possible ordering of joins
 - ➔ Determine the cost of each ordering
 - ➔ Choose the join ordering with minimal cost

Static Algorithm

For joins, two alternative algorithms :

- Nested loops

```
for each tuple of external relation (cardinality  $n_1$ )  
  for each tuple of internal relation (cardinality  $n_2$ )  
    join two tuples if the join predicate is true  
  end  
end
```

→ Complexity: $n_1 * n_2$

- Merge join

```
sort relations  
merge relations
```

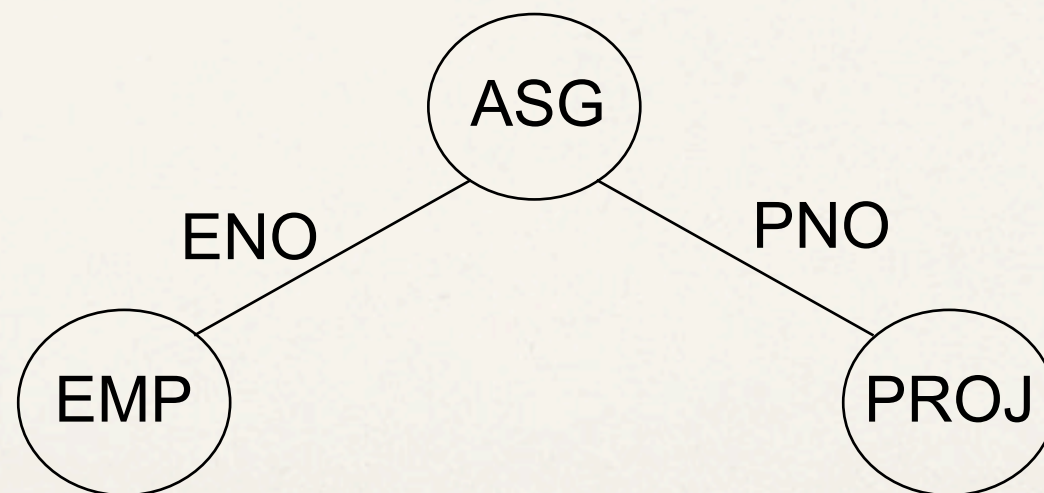
→ Complexity: $n_1 + n_2$ if relations are previously sorted and equijoin

Static Algorithm – Example

Names of employees working on the CAD/CAM project

Assume

- ➔ EMP has an index on ENO,
- ➔ ASG has an index on PNO,
- ➔ PROJ has an index on PNO and an index on PNAME

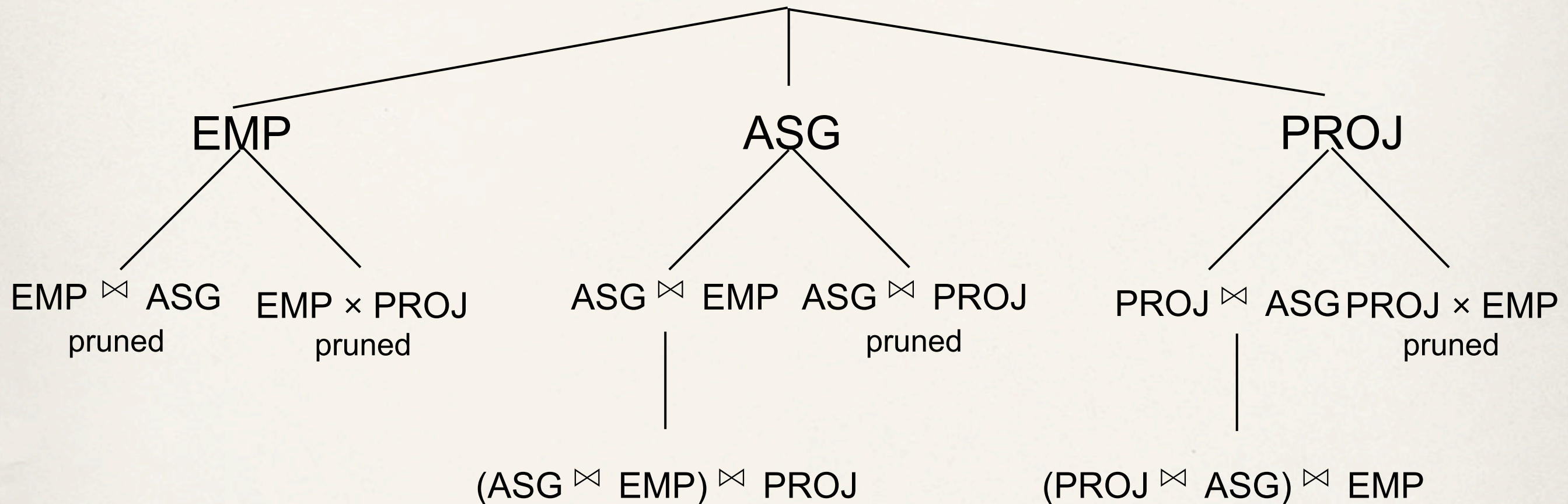


Example (cont'd)

- ① Choose the best access paths to each relation
 - EMP: sequential scan (no selection on EMP)
 - ASG: sequential scan (no selection on ASG)
 - PROJ: index on PNAME (there is a selection on PROJ based on PNAME)
- ② Determine the best join ordering
 - $EMP \bowtie ASG \bowtie PROJ$
 - $ASG \bowtie PROJ \bowtie EMP$
 - $PROJ \bowtie ASG \bowtie EMP$
 - $ASG \bowtie EMP \bowtie PROJ$
 - $EMP \times PROJ \bowtie ASG$
 - $PRO \times JEMP \bowtie ASG$
 - Select the best ordering based on the join costs evaluated according to the two methods

Static Algorithm

Alternatives



Best total join order is one of

((ASG \bowtie EMP) \bowtie PROJ)

((PROJ \bowtie ASG) \bowtie EMP)

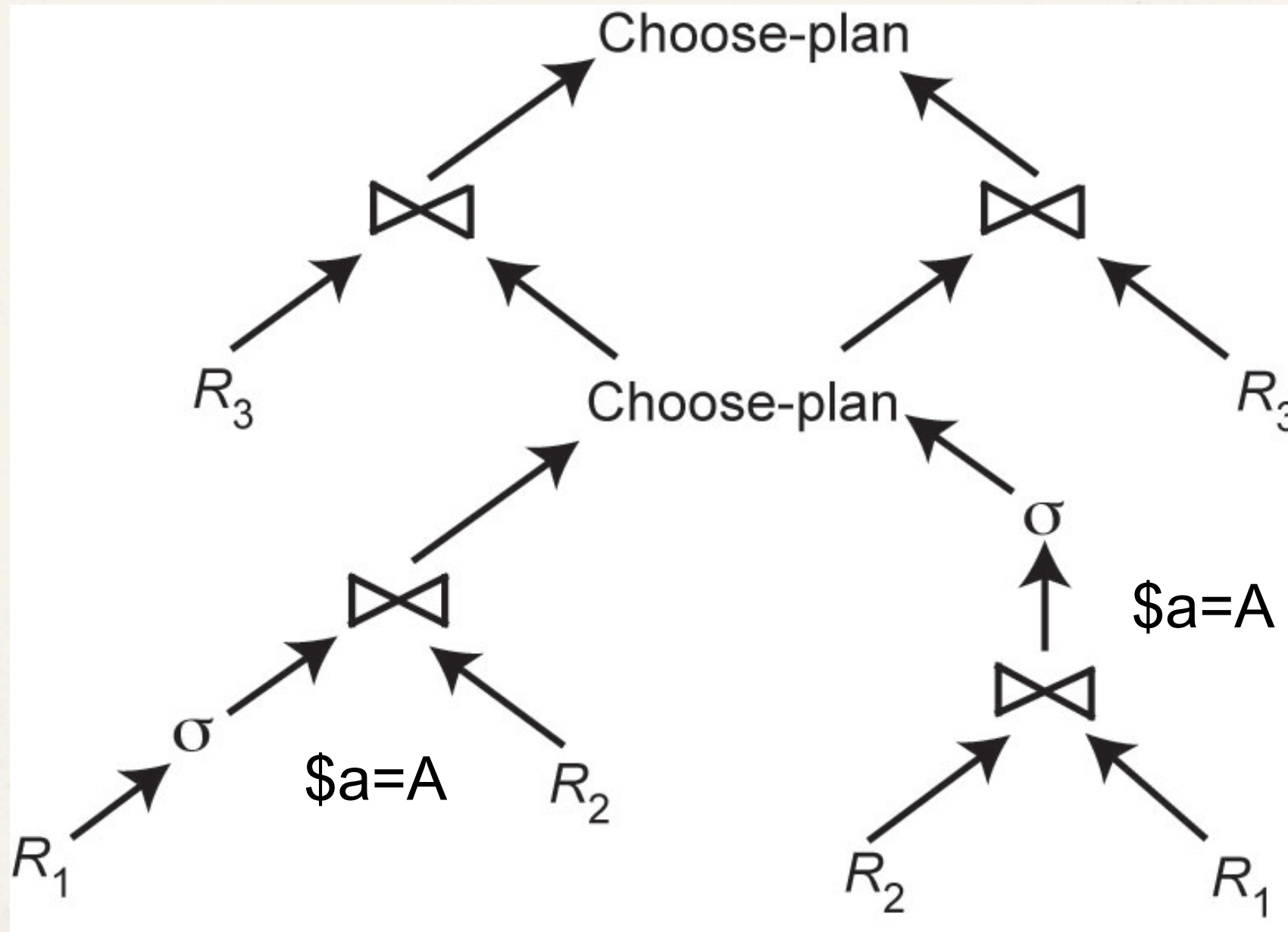
Static Algorithm

- $((\text{PROJ} \bowtie \text{ASG}) \bowtie \text{EMP})$ has a useful index on the select attribute and direct access to the join attributes of ASG and EMP
- Therefore, chose it with the following access methods:
 - ➔ select PROJ using index on PNAME
 - ➔ then join with ASG using index on PNO
 - ➔ then join with EMP using index on ENO

Hybrid optimization

- In general, static optimization is more efficient than dynamic optimization
 - ➔ Adopted by all commercial DBMS
- But even with a sophisticated cost model (with histograms), accurate cost prediction is difficult
- Example
 - ➔ Consider a parametric query with predicate
WHERE R.A = \$a /* \$a is a parameter
 - ➔ The only possible assumption at compile time is uniform distribution of values
- Solution: Hybrid optimization
 - ➔ Choose-plan done at runtime, based on the actual parameter binding

Hybrid Optimization Example

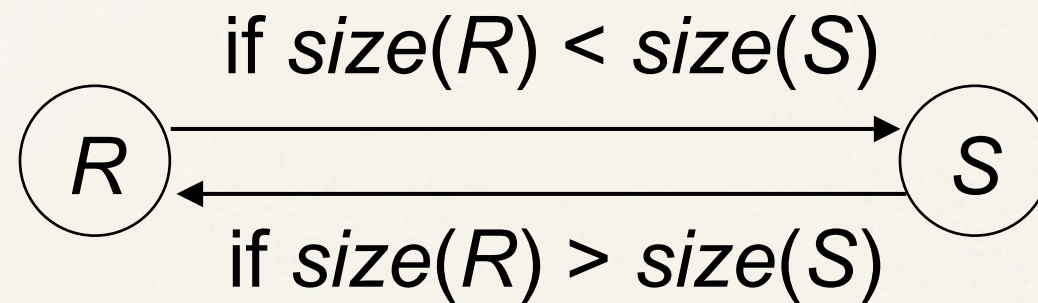


Join Ordering in Fragment Queries

- Ordering joins
 - ➔ Distributed INGRES
 - ➔ System R*
 - ➔ Two-step
- Semijoin ordering
 - ➔ SDD-1

Join Ordering

- Consider two relations only

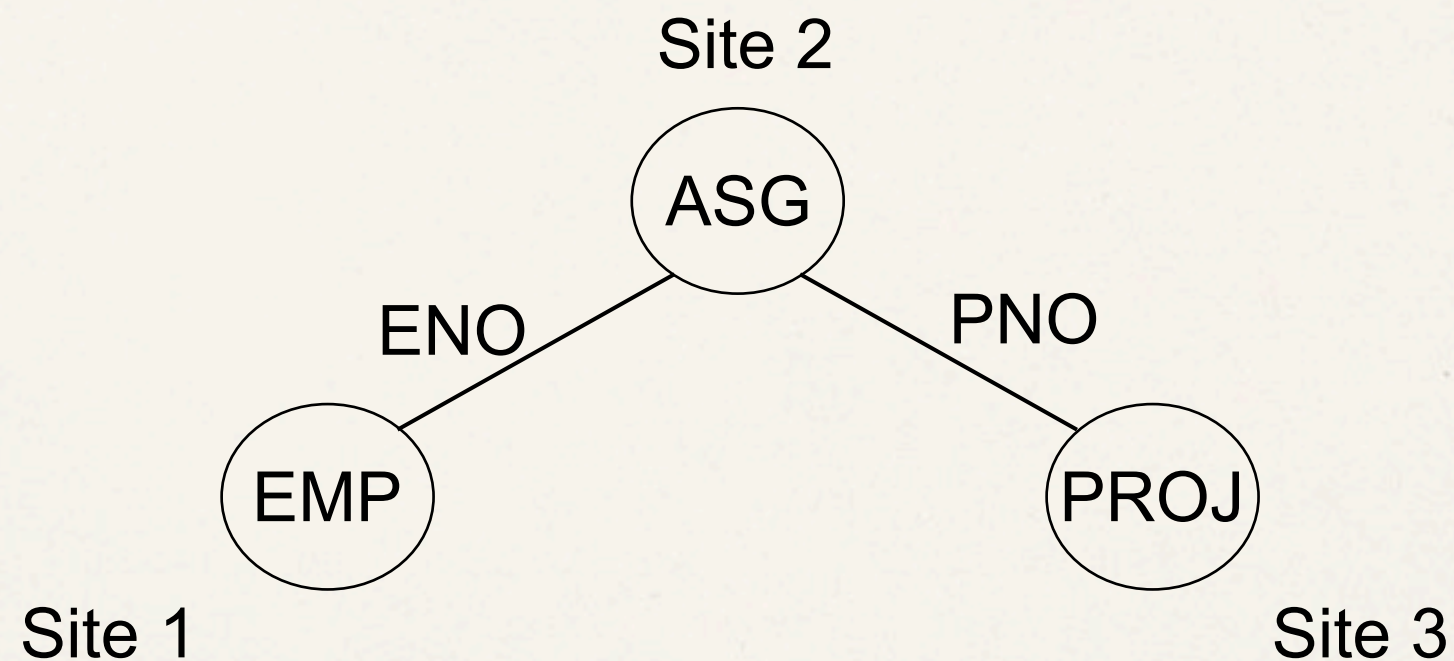


- Multiple relations more difficult because too many alternatives.
 - ➔ Compute the cost of all alternatives and select the best one.
 - ♦ Necessary to compute the size of intermediate relations which is difficult.
 - ➔ Use heuristics

Join Ordering – Example

Consider

$\text{PROJ} \bowtie_{\text{PNO}} \text{ASG} \bowtie_{\text{ENO}} \text{EMP}$



Join Ordering – Example

Execution alternatives:

1. EMP \rightarrow Site 2

Site 2 computes $EMP' = EMP \bowtie ASG$

$EMP' \rightarrow$ Site 3

Site 3 computes $EMP' \bowtie PROJ$

2. ASG \rightarrow Site 1

Site 1 computes $EMP' = EMP \bowtie ASG$

$EMP' \rightarrow$ Site 3

Site 3 computes $EMP' \bowtie PROJ$

3. ASG \rightarrow Site 3

Site 3 computes $ASG' = ASG \bowtie PROJ$

$ASG' \rightarrow$ Site 1

Site 1 computes $ASG' \bowtie EMP$

4. PROJ \rightarrow Site 2

Site 2 computes $PROJ' = PROJ \bowtie ASG$

$PROJ' \rightarrow$ Site 1

Site 1 computes $PROJ' \bowtie EMP$

5. EMP \rightarrow Site 2

PROJ \rightarrow Site 2

Site 2 computes $EMP \bowtie PROJ \bowtie ASG$

Semijoin Algorithms

- Consider the join of two relations:
 - ➔ $R[A]$ (located at site 1)
 - ➔ $S[A]$ (located at site 2)
- Alternatives:
 1. Do the join $R \bowtie_A S$
 2. Perform one of the semijoin equivalents

$$\begin{aligned} R \bowtie_A S &\Leftrightarrow (R \ltimes_A S) \bowtie_A S \\ &\Leftrightarrow R \bowtie_A (S \ltimes_A R) \\ &\Leftrightarrow (R \ltimes_A S) \bowtie_A (S \ltimes_A R) \end{aligned}$$

Semijoin Algorithms

- Perform the join
 - send R to Site 2
 - Site 2 computes $R \bowtie_A S$
- Consider semijoin $(R \bowtie_A S) \bowtie_A S$
 - $S' = \Pi_A(S)$
 - $S' \rightarrow$ Site 1
 - Site 1 computes $R' = R \bowtie_A S'$
 - $R' \rightarrow$ Site 2
 - Site 2 computes $R' \bowtie_A S$

Semijoin is better if

$$size(\Pi_A(S)) + size(R \bowtie_A S) < size(R)$$

Distributed Dynamic Algorithm

1. Execute all monorelation queries (e.g., selection, projection)
2. Reduce the multirelation query to produce irreducible subqueries
 $q_1 \rightarrow q_2 \rightarrow \dots \rightarrow q_n$ such that there is only one relation between q_i and q_{i+1}
1. Choose q_i involving the smallest fragments to execute (call MRQ')
2. Find the best execution strategy for MRQ'
 - a) Determine processing site
 - b) Determine fragments to move
3. Repeat 3 and 4

Static Approach

- Cost function includes local processing as well as transmission
- Considers only joins
- “Exhaustive” search
- Compilation
- Published papers provide solutions to handling horizontal and vertical fragmentations but the implemented prototype does not

Static Approach – Performing Joins

- Ship whole
 - ➔ Larger data transfer
 - ➔ Smaller number of messages
 - ➔ Better if relations are small
- Fetch as needed
 - ➔ Number of messages = $O(\text{cardinality of external relation})$
 - ➔ Data transfer per message is minimal
 - ➔ Better if relations are large and the selectivity is good

Static Approach – Vertical Partitioning & Joins

1. Move outer relation tuples to the site of the inner relation
 - (a) Retrieve outer tuples
 - (b) Send them to the inner relation site
 - (c) Join them as they arrive

$$\begin{aligned}\text{Total Cost} = & \text{cost}(\text{retrieving qualified outer tuples}) \\ & + \text{no. of outer tuples fetched} * \text{cost}(\text{retrieving qualified} \\ & \quad \text{inner tuples}) \\ & + \text{msg. cost} * (\text{no. outer tuples fetched} * \text{avg. outer} \\ & \quad \text{tuple size}) / \text{msg. size}\end{aligned}$$

Static Approach – Vertical Partitioning & Joins

2. Move inner relation to the site of outer relation

Cannot join as they arrive; they need to be stored

$$\begin{aligned} \text{Total cost} = & \text{cost}(\text{retrieving qualified outer tuples}) \\ & + \text{no. of outer tuples fetched} * \text{cost}(\text{retrieving} \\ & \quad \text{matching inner tuples from temporary storage}) \\ & + \text{cost}(\text{retrieving qualified inner tuples}) \\ & + \text{cost}(\text{storing all qualified inner tuples in temporary} \\ & \quad \text{storage}) \\ & + \text{msg. cost} * \text{no. of inner tuples fetched} * \text{avg. inner} \\ & \quad \text{tuple size/msg. size} \end{aligned}$$

Static Approach – Vertical Partitioning & Joins

3. Move both inner and outer relations to another site

$$\begin{aligned}\text{Total cost} = & \text{cost}(\text{retrieving qualified outer tuples}) \\ & + \text{cost}(\text{retrieving qualified inner tuples}) \\ & + \text{cost}(\text{storing inner tuples in storage}) \\ & + \text{msg. cost} \cdot (\text{no. of outer tuples fetched} * \text{avg. outer} \\ & \quad \text{tuple size}) / \text{msg. size} \\ & + \text{msg. cost} * (\text{no. of inner tuples fetched} * \text{avg. inner} \\ & \quad \text{tuple size}) / \text{msg. size} \\ & + \text{no. of outer tuples fetched} * \text{cost}(\text{retrieving inner} \\ & \quad \text{tuples from temporary storage})\end{aligned}$$

Static Approach – Vertical Partitioning & Joins

4. Fetch inner tuples as needed

- (a) Retrieve qualified tuples at outer relation site
- (b) Send request containing join column value(s) for outer tuples to inner relation site
- (c) Retrieve matching inner tuples at inner relation site
- (d) Send the matching inner tuples to outer relation site
- (e) Join as they arrive

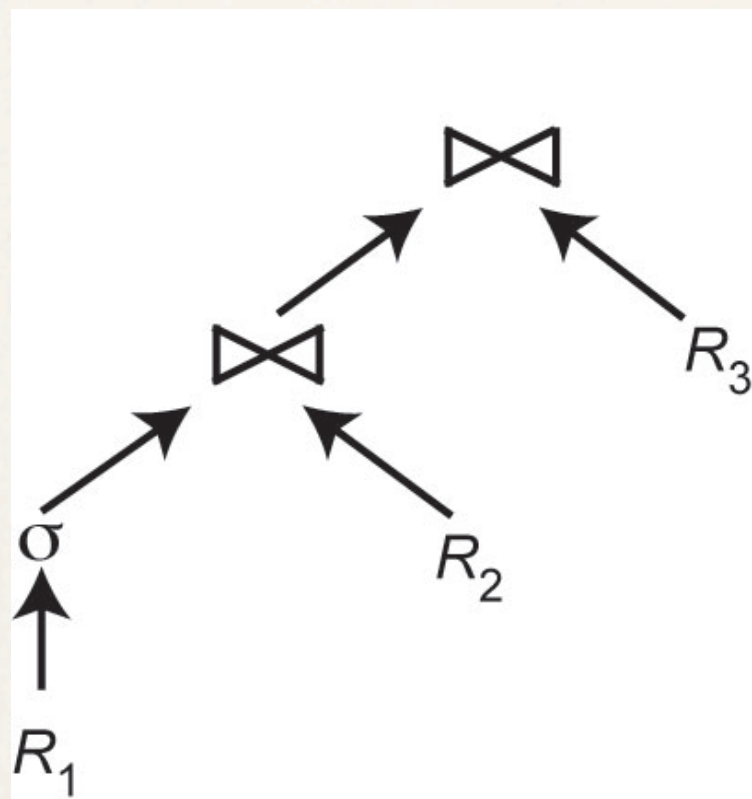
$$\begin{aligned} \text{Total Cost} = & \text{cost}(\text{retrieving qualified outer tuples}) \\ & + \text{msg. cost} * (\text{no. of outer tuples fetched}) \\ & + \text{no. of outer tuples fetched} * \text{no. of} \\ & \quad \text{inner tuples fetched} * \text{avg. inner tuple} \\ & \quad \text{size} * \text{msg. cost} / \text{msg. size}) \\ & + \text{no. of outer tuples fetched} * \text{cost}(\text{retrieving} \\ & \quad \text{matching inner tuples for one outer value}) \end{aligned}$$

Dynamic vs. Static vs Semijoin

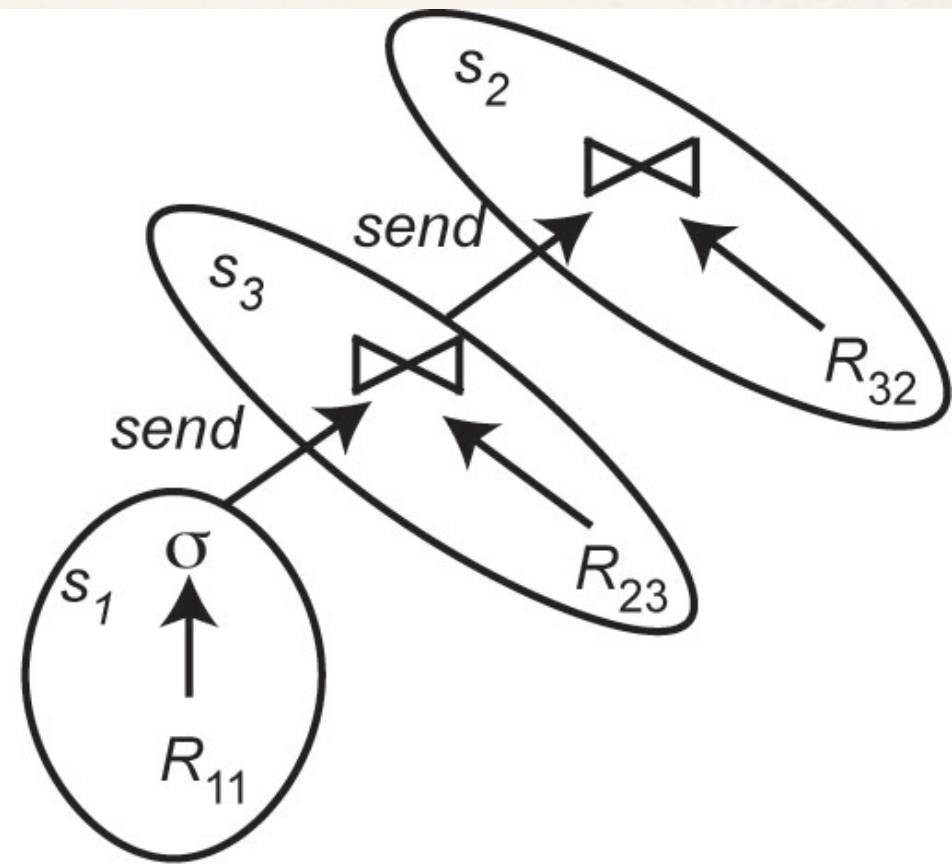
- Semijoin
 - ➔ SDD1 selects only locally optimal schedules
- Dynamic and static approaches have the same advantages and drawbacks as in centralized case
 - ➔ But the problems of accurate cost estimation at compile-time are more severe
 - ◆ More variations at runtime
 - ◆ Relations may be replicated, making site and copy selection important
- Hybrid optimization
 - ➔ Choose-plan approach can be used
 - ➔ 2-step approach simpler

2-Step Optimization

1. At compile time, generate a static plan with operation ordering and access methods only
2. At startup time, carry out site and copy selection and allocate operations to sites



(a) Static plan



(b) Run-time plan

2-Step – Problem Definition

- Given

- ➔ A set of sites $S = \{s_1, s_2, \dots, s_n\}$ with the load of each site
 - ➔ A query $Q = \{q_1, q_2, q_3, q_4\}$ such that each subquery q_i is the maximum processing unit that accesses one relation and communicates with its neighboring queries
 - ➔ For each q_i in Q , a feasible allocation set of sites $S_q = \{s_1, s_2, \dots, s_k\}$ where each site stores a copy of the relation in q_i
- The objective is to find an optimal allocation of Q to S such that
 - ➔ the load unbalance of S is minimized
 - ➔ The total communication cost is minimized

2-Step Algorithm

- For each q in Q compute load (S_q)
- While Q not empty do
 1. Select subquery a with least allocation flexibility
 2. Select best site b for a (with least load and best benefit)
 3. Remove a from Q and recompute loads if needed

2-Step Algorithm Example

- Let $Q = \{q_1, q_2, q_3, q_4\}$ where q_1 is associated with R_1 , q_2 is associated with R_2 joined with the result of q_1 , etc.
- Iteration 1: select q_4 , allocate to s_1 , set $\text{load}(s_1)=2$
- Iteration 2: select q_2 , allocate to s_2 , set $\text{load}(s_2)=3$
- Iteration 3: select q_3 , allocate to s_1 , set $\text{load}(s_1) = 3$
- Iteration 4: select q_1 , allocate to s_3 or s_4

sites	load	R_1	R_2	R_3	R_4
s_1	1	R_{11}		R_{31}	R_{41}
s_2	2		R_{22}		
s_3	2	R_{13}		R_{33}	
s_4	2	R_{14}	R_{24}		

Note: if in iteration 2, q_2 , were allocated to s_4 , this would have produced a better plan. So hybrid optimization can still miss optimal plans