

Distributed Database Management Systems

Outline

- Introduction
- Distributed DBMS Architecture
- Distributed Database Design
- Distributed Query Processing
- Distributed Concurrency Control
- Distributed Reliability Protocols

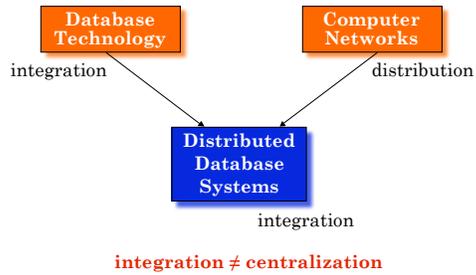
Distributed DBMS

Outline

- Introduction
 - What is a distributed DBMS
 - Problems
 - Current state-of-affairs
- Distributed DBMS Architecture
- Distributed Database Design
- Distributed Query Processing
- Distributed Concurrency Control
- Distributed Reliability Protocols

Distributed DBMS

Motivation



Distributed DBMS

What is a Distributed Database System?

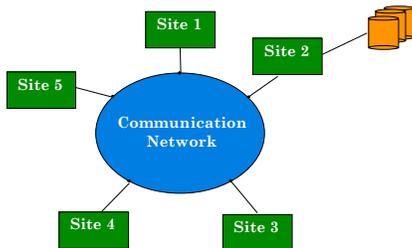
A distributed database (DDB) is a collection of multiple, *logically interrelated* databases distributed over a *computer network*.

A distributed database management system (D-DBMS) is the software that manages the DDB and provides an access mechanism that makes this distribution *transparent* to the users.

Distributed database system (DDBS) = DDB + D-DBMS

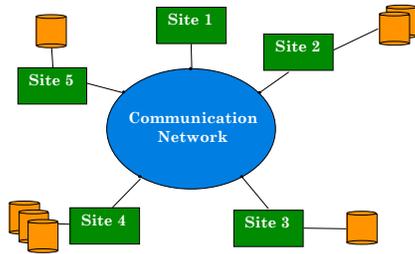
Distributed DBMS

Centralized DBMS on Network



Distributed DBMS

Distributed DBMS Environment



Distributed DBMS

Implicit Assumptions

- Data stored at a number of sites \Rightarrow each site *logically* consists of a single processor.
- Processors at different sites are *interconnected by a computer network* \Rightarrow no multiprocessors
 - parallel database systems
- Distributed database is a *database, not a collection of files* \Rightarrow data logically related as exhibited in the users' access patterns
 - relational data model
- D-DBMS is a *full-fledged DBMS*
 - not remote file system, not a TP system

Distributed DBMS

Distributed DBMS Promises

- 1 Transparent management of distributed, fragmented, and replicated data
- 2 Improved reliability/availability through distributed transactions
- 3 Improved performance
- 4 Easier and more economical system expansion

Distributed DBMS

Transparency

- Transparency is the separation of the higher level semantics of a system from the lower level implementation issues.
- Fundamental issue is to provide **data independence** in the distributed environment
 - Network (distribution) transparency
 - Replication transparency
 - Fragmentation transparency
 - ◆ horizontal fragmentation: selection
 - ◆ vertical fragmentation: projection
 - ◆ hybrid

Distributed DBMS

Example

EMP			ASG			
ENO	ENAME	TITLE	ENO	PNO	RESP	DUR
E1	J. Doe	Elect. Eng.	E1	P1	Manager	12
E2	M. Smith	Syst. Anal.	E2	P1	Analyst	24
E3	A. Lee	Mech. Eng.	E2	P2	Analyst	6
E4	J. Miller	Programmer	E3	P3	Consultant	10
E5	B. Casey	Syst. Anal.	E3	P4	Engineer	48
E6	L. Chu	Elect. Eng.	E4	P2	Programmer	18
E7	R. Davis	Mech. Eng.	E5	P2	Manager	24
E8	J. Jones	Syst. Anal.	E6	P4	Manager	48
			E7	P3	Engineer	36
			E7	P5	Engineer	23
			E8	P3	Manager	40

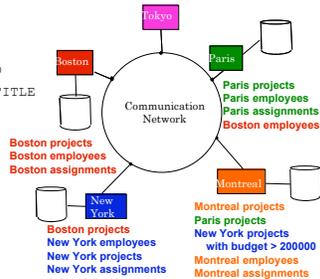
PROJ			PAY	
PNO	PNAME	BUDGET	TITLE	SAL
P1	Instrumentation	150000	Elect. Eng.	40000
P2	Database Develop.	135000	Syst. Anal.	34000
P3	CAD/CAM	250000	Mech. Eng.	27000
P4	Maintenance	310000	Programmer	24000

Distributed DBMS

Transparent Access

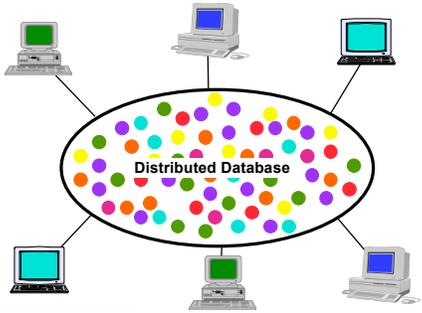
```

SELECT ENAME, SAL
FROM EMP, ASG, PAY
WHERE DUR > 12
AND EMP.ENO = ASG.ENO
AND PAY.TITLE = EMP.TITLE
    
```



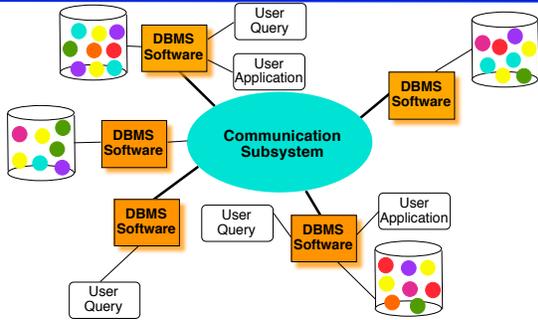
Distributed DBMS

Distributed Database – User View



Distributed DBMS

Distributed DBMS - Reality



Distributed DBMS

Potentially Improved Performance

- Proximity of data to its points of use
 - Requires some support for fragmentation and replication
- Parallelism in execution
 - Inter-query parallelism
 - Intra-query parallelism

Distributed DBMS

Parallelism Requirements

- Have as much of the data required by *each* application at the site where the application executes
 - Full replication
- How about updates?
 - Updates to replicated data requires implementation of distributed concurrency control and commit protocols

Distributed DBMS

System Expansion

- Issue is database scaling
- Emergence of microprocessor and workstation technologies
 - Demise of Grosh's law
 - Client-server model of computing
- Data communication cost vs telecommunication cost

Distributed DBMS

Distributed DBMS Issues

- **Distributed Database Design**
 - how to distribute the database
 - replicated & non-replicated database distribution
 - a related problem in directory management
- **Query Processing**
 - convert user transactions to data manipulation instructions
 - optimization problem
 - $\min\{\text{cost} = \text{data transmission} + \text{local processing}\}$
 - general formulation is NP-hard

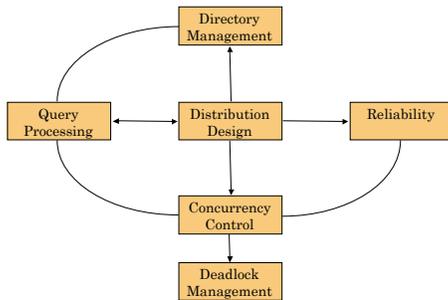
Distributed DBMS

Distributed DBMS Issues

- **Concurrency Control**
 - synchronization of concurrent accesses
 - consistency and isolation of transactions' effects
 - deadlock management
- **Reliability**
 - how to make the system resilient to failures
 - atomicity and durability

Distributed DBMS

Relationship Between Issues



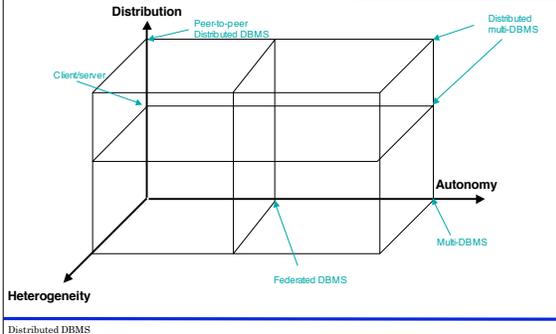
Distributed DBMS

Outline

- Introduction
- Distributed DBMS Architecture
 - Implementation Alternatives
 - Component Architecture
- Distributed Database Design
- Distributed Query Processing
- Distributed Concurrency Control
- Distributed Reliability Protocols

Distributed DBMS

DBMS Implementation Alternatives

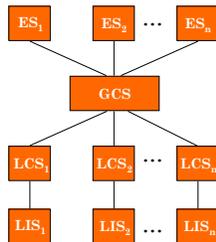


Dimensions of the Problem

- **Distribution**
 - Whether the components of the system are located on the same machine or not
- **Heterogeneity**
 - Various levels (hardware, communications, operating system)
 - DBMS important one
 - data model, query language, transaction management algorithms
- **Autonomy**
 - Not well understood and most troublesome
 - Various versions
 - **Design autonomy:** Ability of a component DBMS to decide on issues related to its own design.
 - **Communication autonomy:** Ability of a component DBMS to decide whether and how to communicate with other DBMSs.
 - **Execution autonomy:** Ability of a component DBMS to execute local operations in any manner it wants to.

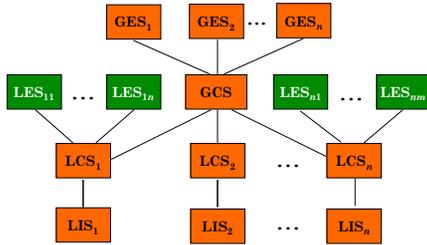
Distributed DBMS

Datological Distributed DBMS Architecture



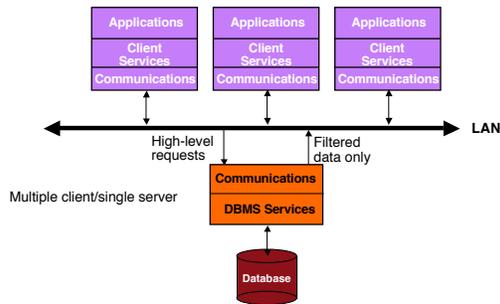
Distributed DBMS

Datalogical Multi-DBMS Architecture



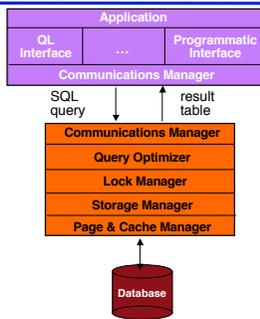
Distributed DBMS

Clients/Server



Distributed DBMS

Task Distribution



Distributed DBMS

Advantages of Client-Server Architectures

- More efficient division of labor
- Horizontal and vertical scaling of resources
- Better price/performance on client machines
- Ability to use familiar tools on client machines
- Client access to remote data (via standards)
- Full DBMS functionality provided to client workstations
- Overall better system price/performance

Distributed DBMS

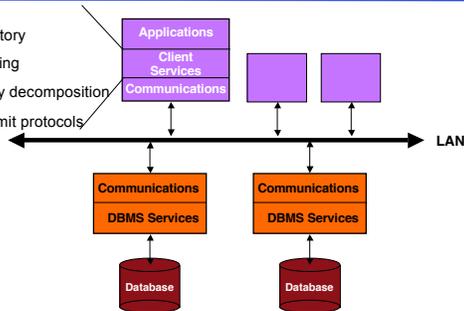
Problems With Multiple-Client/Single Server

- Server forms bottleneck
- Server forms single point of failure
- Database scaling difficult

Distributed DBMS

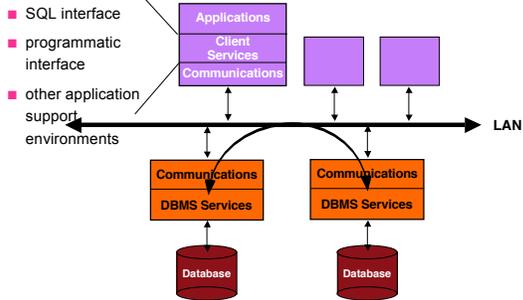
Multiple Clients/Multiple Servers

- directory
- caching
- query decomposition
- commit protocols



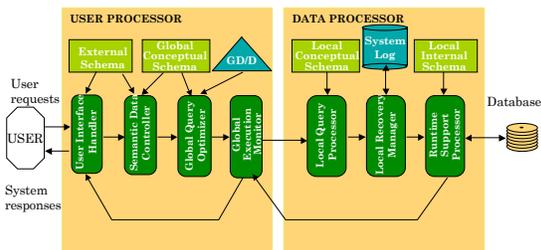
Distributed DBMS

Server-to-Server



Distributed DBMS

Peer-to-Peer Component Architecture



Distributed DBMS

Outline

- Introduction
- Distributed DBMS Architecture
- Distributed Database Design
 - Fragmentation
 - Data Placement
- Distributed Query Processing
- Distributed Concurrency Control
- Distributed Reliability Protocols

Distributed DBMS

Design Problem

- In the general setting :
 - Making decisions about the placement of data and programs across the sites of a computer network as well as possibly designing the network itself.
- In Distributed DBMS, the placement of applications entails
 - placement of the distributed DBMS software; and
 - placement of the applications that run on the database

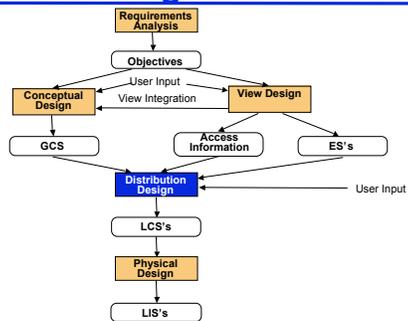
Distributed DBMS

Distribution Design

- Top-down
 - mostly in designing systems from scratch
 - mostly in homogeneous systems
- Bottom-up
 - when the databases already exist at a number of sites

Distributed DBMS

Top-Down Design



Distributed DBMS

Distribution Design

- Fragmentation
 - Localize access
 - Horizontal fragmentation
 - Vertical fragmentation
 - Hybrid fragmentation
- Distribution
 - Placement of fragments on nodes of a network

Distributed DBMS

Horizontal Fragmentation

PROJ₁ : projects with budgets less than \$200,000
 PROJ₂ : projects with budgets greater than or equal to \$200,000

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris
P5	CAD/CAM	500000	Boston

PROJ ₁			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York

PROJ ₂			
PNO	PNAME	BUDGET	LOC
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris
P5	CAD/CAM	500000	Boston

Distributed DBMS

Vertical Fragmentation

PROJ₁: information about project budgets

PROJ₂: information about project names and locations

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris
P5	CAD/CAM	500000	Boston

PROJ ₁	
PNO	BUDGET
P1	150000
P2	135000
P3	250000
P4	310000
P5	500000

PROJ ₂		
PNO	PNAME	LOC
P1	Instrumentation	Montreal
P2	Database Develop.	New York
P3	CAD/CAM	New York
P4	Maintenance	Paris
P5	CAD/CAM	Boston

Distributed DBMS

Correctness of Fragmentation

- **Completeness**
 - Decomposition of relation R into fragments R_1, R_2, \dots, R_n is complete iff each data item in R can also be found in some R_i
- **Reconstruction**
 - If relation R is decomposed into fragments R_1, R_2, \dots, R_n , then there should exist some relational operator ∇ such that
$$R = \nabla_{1 \leq i \leq n} R_i$$
- **Disjointness**
 - If relation R is decomposed into fragments R_1, R_2, \dots, R_n , and data item d_i is in R_j , then d_i should not be in any other fragment R_k ($k \neq j$).

Distributed DBMS

Allocation Alternatives

- **Non-replicated**
 - partitioned : each fragment resides at only one site
- **Replicated**
 - fully replicated : each fragment at each site
 - partially replicated : each fragment at some of the sites
- **Rule of thumb:**

If $\frac{\text{read - only queries}}{\text{update queries}} \geq 1$ replication is advantageous,
otherwise replication may cause problems

Distributed DBMS

Fragment Allocation

- **Problem Statement**
 - Given
 - ◆ $F = \{F_1, F_2, \dots, F_n\}$ fragments
 - ◆ $S = \{S_1, S_2, \dots, S_m\}$ network sites
 - ◆ $Q = \{q_1, q_2, \dots, q_q\}$ applications
 - Find the "optimal" distribution of F to S .
- **Optimality**
 - Minimal cost
 - ◆ Communication + storage + processing (read & update)
 - ◆ Cost in terms of time (usually)
 - Performance
 - ◆ Response time and/or throughput
 - Constraints
 - ◆ Per site constraints (storage & processing)

Distributed DBMS

Allocation Model

General Form

min(Total Cost)
subject to
response time constraint
storage constraint
processing constraint

Decision Variable

$x_{ij} = \begin{cases} 1 & \text{if fragment } F_i \text{ is stored at site } S_j \\ 0 & \text{otherwise} \end{cases}$

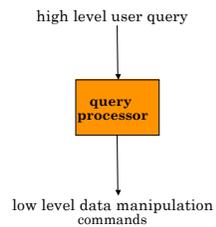
Distributed DBMS

Outline

- Introduction
- Distributed DBMS Architecture
- Distributed Database Design
- Distributed Query Processing
 - Query Processing Methodology
 - Distributed Query Optimization
- Distributed Concurrency Control
- Distributed Reliability Protocols

Distributed DBMS

Query Processing



Distributed DBMS

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?

Distributed DBMS

Selecting Alternatives

```

SELECT      ENAME
FROM        EMP, ASG
WHERE       EMP.ENO = ASG.ENO
AND        DUR > 37
    
```

Strategy 1

$$\Pi_{ENAME}(\sigma_{DUR>37 \wedge EMP.ENO=ASG.ENO}(EMP \times ASG))$$

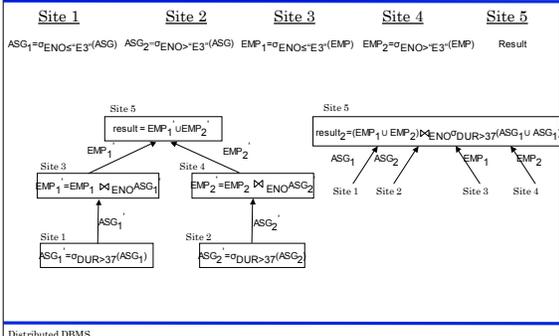
Strategy 2

$$\Pi_{ENAME}(EMP \bowtie_{ENO} (\sigma_{DUR>37}(ASG)))$$

Strategy 2 avoids Cartesian product, so is "better"

Distributed DBMS

What is the Problem?



Distributed DBMS

Cost of Alternatives

Assume:

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

Strategy 1

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

Strategy 2

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

Distributed DBMS

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 will dominate
 - low bandwidth
 - low speed
 - high protocol overhead
- most algorithms ignore all other cost components

Local area networks

- communication cost not that dominant
- total cost function should be considered

Can also maximize throughput

Distributed DBMS

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

Distributed DBMS

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

Distributed DBMS

Query Optimization Issues – Optimization Timing

- Static
 - compilation ⇒ optimize prior to the execution
 - difficult to estimate the size of the intermediate results ⇒ 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

Distributed DBMS

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

Distributed DBMS

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

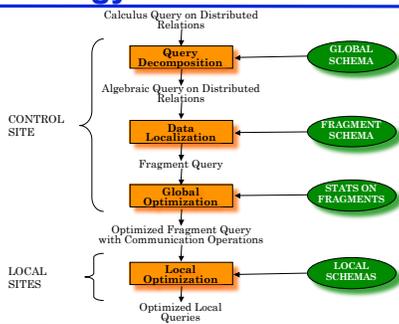
Distributed DBMS

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 DBMS

Distributed Query Processing Methodology



Distributed DBMS

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 \Rightarrow algebraic query
 - more than one translation is possible
 - use transformation rules

Distributed DBMS

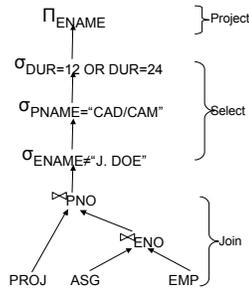
Restructuring

- 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  $\neq$  "J. Doe"
AND PNAME = "CAD/CAM"
AND (DUR = 12 OR DUR = 24)
    
```



Distributed DBMS

Restructuring – Transformation Rules (Examples)

- **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**
 - $\Pi_{A'}(\Pi_{A'}(R)) \Leftrightarrow \Pi_{A'}(R)$
 - $\sigma_{p_1(A_1)}(\sigma_{p_2(A_2)}(R)) = \sigma_{p_1(A_1) \wedge p_2(A_2)}(R)$
where $R[A]$ and $A' \subseteq A$, $A'' \subseteq A$ and $A' \subseteq A''$
- **Commuting selection with projection**

Distributed DBMS

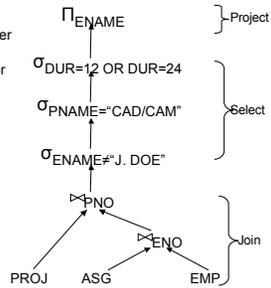
Example

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.

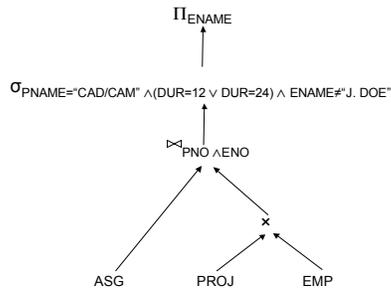
```

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)
    
```



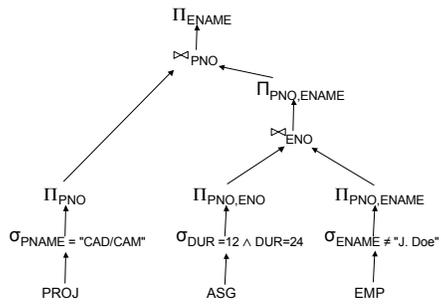
Distributed DBMS

Equivalent Query



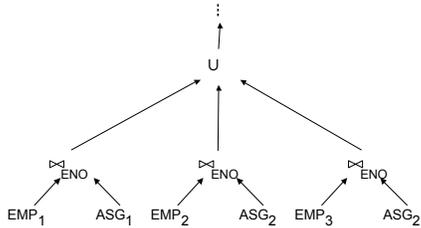
Distributed DBMS

Restructuring



Distributed DBMS

Eliminates Unnecessary Work



Distributed DBMS

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)

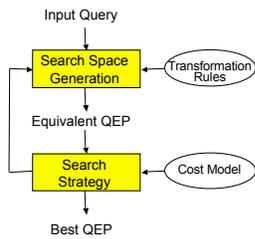
Distributed DBMS

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,...)

Distributed DBMS

Query Optimization Process



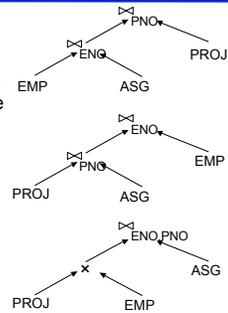
Distributed DBMS

Search Space

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

```

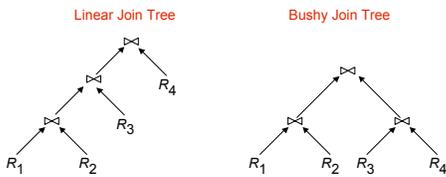
SELECT ENAME, RESP
FROM EMP, ASG, PROJ
WHERE EMP.ENO=ASG.ENO
AND ASG.PNO=PROJ.PNO
  
```



Distributed DBMS

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



Distributed DBMS

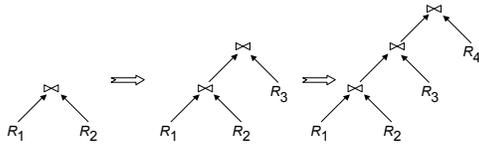
Search Strategy

- 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 > 5-6 relations
 - Simulated annealing
 - Iterative improvement

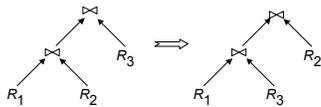
Distributed DBMS

Search Strategies

- Deterministic



- Randomized



Distributed DBMS

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

Distributed DBMS

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

Distributed DBMS

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

Distributed DBMS

Response Time

Elapsed time between the initiation and the completion of a query

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

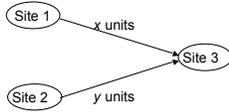
CPU time = unit instruction time * no. of sequential instructions

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

Distributed DBMS

Example

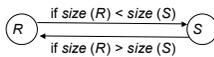


Assume that only the communication cost is considered
 Total time = 2 * message initialization time + unit transmission time * (x+y)
 Response time = max {time to send x from 1 to 3, time to send y from 2 to 3}
 time to send x from 1 to 3 = message initialization time + unit transmission time * x
 time to send y from 2 to 3 = message initialization time + unit transmission time * y

Distributed DBMS

Join Ordering

- Alternatives
 - Ordering joins
 - Semijoin ordering
- Consider two relations only



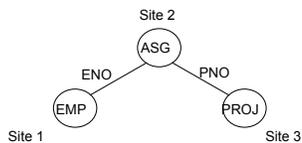
- Multiple relations more difficult because too many alternatives.
 - Compute the cost of all alternatives and select the best one.
 - ◆ Necessary to compute the size of intermediate relations which is difficult.
 - Use heuristics

Distributed DBMS

Join Ordering – Example

Consider

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



Distributed DBMS

Join Ordering – Example

Execution alternatives:

- | | |
|-----------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------|
| 1. EMP → Site 2
Site 2 computes EMP=EMP<A>ASG
EMP' → Site 3
Site 3 computes EMP' >A> PROJ | 2. ASG → Site 1
Site 1 computes EMP=EMP>A> ASG
EMP' → Site 3
Site 3 computes EMP'>A> PROJ |
| 3. ASG → Site 3
Site 3 computes ASG'=ASG >A> PROJ
ASG' → Site 1
Site 1 computes ASG>A> EMP | 4. PROJ → Site 2
Site 2 computes PROJ=PROJ<A>ASG
PROJ' → Site 1
Site 1 computes PROJ'>A> EMP |
| 5. EMP → Site 2
PROJ → Site 2
Site 2 computes EMP >A> PROJ >A> ASG | |

Distributed DBMS

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

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

Distributed DBMS

Semijoin Algorithms

■ Perform the join

- send R to Site 2
- Site 2 computes $R \bowtie_A S$

■ Consider semijoin $(R \ltimes_A S) \triangleright_A S$

- $S' \leftarrow \Pi_A(S)$
- $S' \rightarrow$ Site 1
- Site 1 computes $R' = R \bowtie_A S'$
- $R' \rightarrow$ Site 2
- Site 2 computes $R' \triangleright_A S$

Semijoin is better if

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

Distributed DBMS

R* Algorithm

- 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

Distributed DBMS

R* Algorithm

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

Distributed DBMS

R* Algorithm – 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

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

Distributed DBMS

R* Algorithm – Vertical Partitioning & Joins

2. Move inner relation to the site of outer relation

cannot join as they arrive; they need to be stored

Total Cost = cost(retrieving qualified outer tuples)

- + no. of outer tuples fetched * cost(retrieving matching inner tuples from temporary storage)
- + cost(retrieving qualified inner tuples)
- + cost(storing all qualified inner tuples in temporary storage)
- + msg. cost * (no. of inner tuples fetched * avg. inner tuple size) / msg. size

Distributed DBMS

R* Algorithm – Vertical Partitioning & Joins

3. Move both inner and outer relations to another site

Total cost = cost(retrieving qualified outer tuples)

- + cost(retrieving qualified inner tuples)
- + cost(storing inner tuples in storage)
- + msg. cost * (no. of outer tuples fetched * avg. outer tuple size) / msg. size
- + msg. cost * (no. of inner tuples fetched * avg. inner tuple size) / msg. size
- + no. of outer tuples fetched * cost(retrieving inner tuples from temporary storage)

Distributed DBMS

R* Algorithm – 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

Total Cost = cost(retrieving qualified outer tuples)

- + msg. cost * (no. of outer tuples fetched)
- + no. of outer tuples fetched * (no. of inner tuples fetched * avg. inner tuple size * msg. cost / msg. size)
- + no. of outer tuples fetched * cost(retrieving matching inner tuples for one outer value)

Distributed DBMS

Step 4 – Local Optimization

Input: Best global execution schedule

- Select the best **access path**
- Use the centralized optimization techniques

Distributed DBMS

Outline

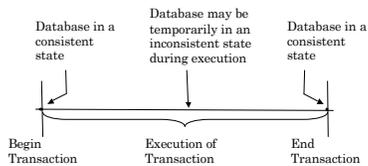
- Introduction
- Distributed DBMS Architecture
- Distributed Database Design
- Distributed Query Processing
- Distributed Concurrency Control
 - Transaction Concepts & Models
 - Serializability
 - Distributed Concurrency Control Protocols
- Distributed Reliability Protocols

Distributed DBMS

Transaction

A transaction is a collection of actions that make consistent transformations of system states while preserving system consistency.

- concurrency transparency
- failure transparency



Distributed DBMS

Example Database

Consider an airline reservation example with the relations:

```
FLIGHT(FNO, DATE, SRC, DEST, STSOLD, CAP)
CUST(CNAME, ADDR, BAL)
FC(FNO, DATE, CNAME, SPECIAL)
```

Distributed DBMS

Example Transaction

```
Begin_transaction Reservation
begin
  input(flight_no, date, customer_name);
  EXEC SQL      UPDATE  FLIGHT
                SET     STSOLD = STSOLD + 1
                WHERE   FNO = flight_no AND DATE = date;
  EXEC SQL      INSERT
                INTO    FC(FNO, DATE, CNAME, SPECIAL);
                VALUES (flight_no, date, customer_name, null);
  output("reservation completed")
end . {Reservation}
```

Distributed DBMS

Termination of Transactions

```
Begin_transaction Reservation
begin
input(flight_no, date, customer_name);
EXEC SQL      SELECT  STSOLD,CAP
                INTO   temp1,temp2
                FROM    FLIGHT
                WHERE   FNO = flight_no AND DATE = date;
if temp1 = temp2 then
  output("no free seats");
Abort
else
EXEC SQL      UPDATE  FLIGHT
                SET     STSOLD = STSOLD + 1
                WHERE   FNO = flight_no AND DATE = date;
EXEC SQL      INSERT
                INTO    FC(FNO, DATE, CNAME, SPECIAL);
                VALUES (flight_no, date, customer_name, null);
  Commit
  output("reservation completed")
endif
end . {Reservation}
```

Distributed DBMS

Properties of Transactions

ATOMICITY

- all or nothing

CONSISTENCY

- no violation of integrity constraints

ISOLATION

- concurrent changes invisible & serializable

DURABILITY

- committed updates persist

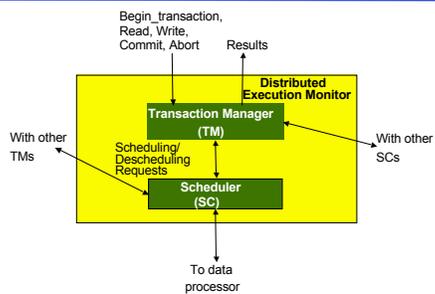
Distributed DBMS

Transactions Provide...

- *Atomic* and *reliable* execution in the presence of failures
- *Correct* execution in the presence of multiple user accesses
- Correct management of *replicas* (if they support it)

Distributed DBMS

Architecture Revisited



Distributed DBMS

Serializable History

- Transactions execute concurrently, but the net effect of the resulting history upon the database is *equivalent* to some *serial* history.
- Equivalent with respect to what?
 - **Conflict equivalence:** the relative order of execution of the conflicting operations belonging to unaborting transactions in two histories are the same.
 - **Conflicting operations:** two incompatible operations (e.g., Read and Write) conflict if they both access the same data item.
 - ◆ Incompatible operations of each transaction is assumed to conflict; do not change their execution orders.
 - ◆ If two operations from two different transactions conflict, the corresponding transactions are also said to conflict.

Distributed DBMS

Serializability in Distributed DBMS

- Somewhat more involved. Two histories have to be considered:
 - local histories
 - global history
- For global transactions (i.e., global history) to be serializable, two conditions are necessary:
 - Each local history should be serializable.
 - Two conflicting operations should be in the same relative order in all of the local histories where they appear together.

Distributed DBMS

Global Non-serializability

T_1 : Read(x)	T_2 : Read(x)
$x \leftarrow x+5$	$x \leftarrow x*15$
Write(x)	Write(x)
Commit	Commit

The following two local histories are individually serializable (in fact serial), but the two transactions are not globally serializable.

$$LH_1 = \{R_1(x), W_1(x), C_1, R_2(x), W_2(x), C_2\}$$

$$LH_2 = \{R_2(x), W_2(x), C_2, R_1(x), W_1(x), C_1\}$$

Distributed DBMS

Concurrency Control Algorithms

- Pessimistic
 - Two-Phase Locking-based (2PL)
 - ◆ Centralized (primary site) 2PL
 - ◆ Primary copy 2PL
 - ◆ Distributed 2PL
 - Timestamp Ordering (TO)
 - ◆ Basic TO
 - ◆ Multiversion TO
 - ◆ Conservative TO
 - Hybrid
- Optimistic
 - Locking-based
 - Timestamp ordering-based

Distributed DBMS

Locking-Based Algorithms

- Transactions indicate their intentions by requesting locks from the scheduler (called **lock manager**).
- Locks are either **read lock (rl)** [also called **shared lock**] or **write lock (wl)** [also called **exclusive lock**]
- Read locks and write locks conflict (because Read and Write operations are incompatible)

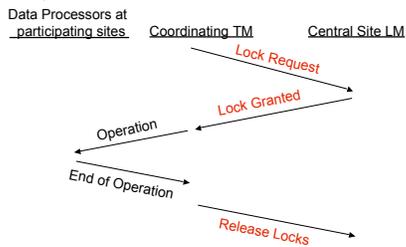
	<i>rl</i>	<i>wl</i>
<i>rl</i>	yes	no
<i>wl</i>	no	no

- Locking works nicely to allow concurrent processing of transactions.

Distributed DBMS

Centralized 2PL

- There is only one 2PL scheduler in the distributed system.
- Lock requests are issued to the central scheduler.



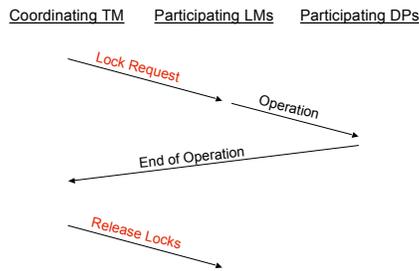
Distributed DBMS

Distributed 2PL

- 2PL schedulers are placed at each site. Each scheduler handles lock requests for data at that site.
- A transaction may read any of the replicated copies of item x , by obtaining a read lock on one of the copies of x . Writing into x requires obtaining write locks for all copies of x .

Distributed DBMS

Distributed 2PL Execution



Distributed DBMS

Timestamp Ordering

- 1 Transaction (T_i) is assigned a globally unique timestamp $ts(T_i)$.
- 2 Transaction manager attaches the timestamp to all operations issued by the transaction.
- 3 Each data item is assigned a write timestamp (wts) and a read timestamp (rts):
 - $rts(x)$ = largest timestamp of any read on x
 - $wts(x)$ = largest timestamp of any read on x
- 4 Conflicting operations are resolved by timestamp order.

Basic T/O:

for $R_i(x)$	for $W_j(x)$
if $ts(T_i) < wts(x)$	if $ts(T_j) < rts(x)$ and $ts(T_j) < wts(x)$
then reject $R_i(x)$	then reject $W_j(x)$
else accept $R_i(x)$	else accept $W_j(x)$
$rts(x) \leftarrow ts(T_i)$	$wts(x) \leftarrow ts(T_j)$

Distributed DBMS

Outline

- Introduction
- Distributed DBMS Architecture
- Distributed Database Design
- Distributed Query Processing
- Distributed Concurrency Control
- Distributed Reliability Protocols
 - Distributed Commit Protocols
 - Distributed Recovery Protocols

Distributed DBMS

Reliability

Problem:

How to maintain

atomicity

durability

properties of transactions

Distributed DBMS

Types of Failures

- Transaction failures
 - Transaction aborts (unilaterally or due to deadlock)
 - Avg. 3% of transactions abort abnormally
- System (site) failures
 - Failure of processor, main memory, power supply, ...
 - Main memory contents are lost, but secondary storage contents are safe
 - Partial vs. total failure
- Media failures
 - Failure of secondary storage devices such that the stored data is lost
 - Head crash/controller failure (?)
- Communication failures
 - Lost/undeliverable messages
 - Network partitioning

Distributed DBMS

Distributed Reliability Protocols

- Commit protocols
 - How to execute commit command for distributed transactions.
 - Issue: how to ensure atomicity and durability?
- Termination protocols
 - If a failure occurs, how can the remaining operational sites deal with it.
 - *Non-blocking* : the occurrence of failures should not force the sites to wait until the failure is repaired to terminate the transaction.
- Recovery protocols
 - When a failure occurs, how do the sites where the failure occurred deal with it.
 - *Independent* : a failed site can determine the outcome of a transaction without having to obtain remote information.
- Independent recovery ⇒ non-blocking termination

Distributed DBMS

Two-Phase Commit (2PC)

Phase 1 : The coordinator gets the participants ready to write the results into the database

Phase 2 : Everybody writes the results into the database

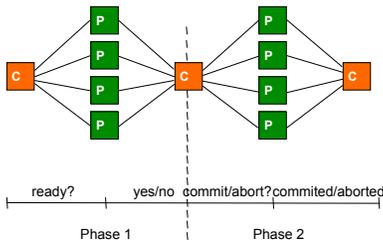
- **Coordinator** : The process at the site where the transaction originates and which controls the execution
- **Participant** : The process at the other sites that participate in executing the transaction

Global Commit Rule:

- ❶ The coordinator aborts a transaction if and only if at least one participant votes to abort it.
- ❷ The coordinator commits a transaction if and only if all of the participants vote to commit it.

Distributed DBMS

Centralized 2PC



Distributed DBMS
