

UNIT 12

Distributed Database

Contents

- ❑ 12.1 Introduction
- ❑ 12.2 The Twelve Objectives
- ❑ 12.3 Problems of Distributed Database Systems
- ❑ 12.4 Gateways
- ❑ 12.5 Client/Server Systems

12.1 Introduction

Distributed Database System

- A system involving multiple sites connected together via communication network.
- User at any site can access data stored at any site.
- Each site is a database system in its own right: its own local database, local users, local DBMS, local DC manager.

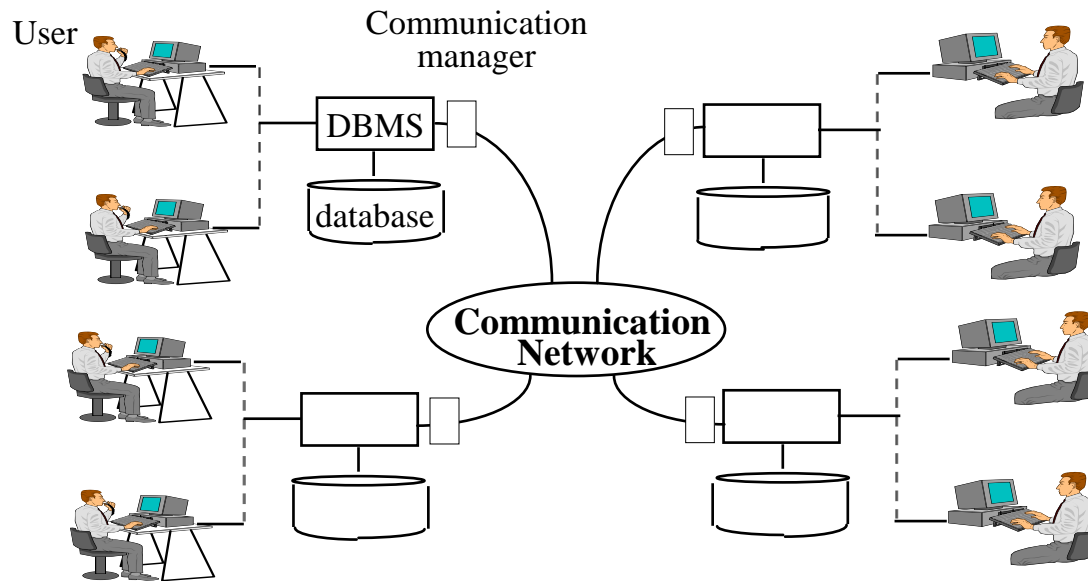


Fig 12.1: A typical distributed database system

Distributed Database System (cont.)

■ Assumptions

- The system is **homogeneous**, in the sense that each site is running its own copy of the same DBMS.
(can be relaxed => **heterogeneous** DBMS)
- The communication network is **slow**: component sites are geographically dispersed.

■ Advantages

- Enables the structure of database to mirror that of the enterprise.
- Combines efficiency of processing with increased accessibility.

■ Disadvantage

- Complex for implementation

Distributed Database System (cont.)

- Sample Systems

- Prototypes

- SDD-1: Computer Corporation of America (CCA), late 1970s and early 1980s.
 - R*: IBM, early 1980s
 - Distributed INGRES: Berkeley, early 1980s.

- Commercial Products :

- INGRES / STAR: Relational Technology Inc.
 - SQL*STAR: Oracle Corp.
 - DB2 version 2 release 2: IBM
 - SQL server: Microsoft

- A Fundamental Principle

To the user, a **distributed system** should look exactly like a **non-distributed system**.

12.2 The Twelve Objectives

The Twelve Objectives

1. Local Autonomy

- all operations at a given site are controlled by that site, should not depend on other sites.
- local data is locally owned and managed.
- Not wholly achievable => sites should be autonomous to the maximum extend possible.

2. No Reliance on a Central Site

- all sites must be treated as equals.
- the central site may be bottleneck.

3. Continuous Operation

- Reliability
- Availability
- Never require the system to be shutdown to perform some function:
e.g. add a new site.

The Twelve Objectives (cont.)

★ 4. Location Independence (Location Transparency)

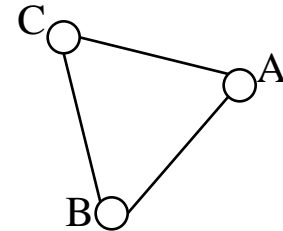
- user should not need to know at which site the data is stored, but should be able to behave as if the entire database were stored at their own local site.
- a request for some remote data => system should find the data automatically.

- Advantages

<1> Simplify user programs and activities

<e.g.>

```
SELECT S#  
FROM S  
AT SITE A  
WHERE SNAME = 'John'
```



<2> allow data to be moved from one site to another at any time without invalidating any program or activities.

The Twelve Objectives (cont.)

★ 5. Fragmentation Independence (Fragmentation Transparency)

- Data Fragmentation
 - a given local object can be divided up into pieces (**fragments**) for physical storage purpose.

<e.g.> user perception

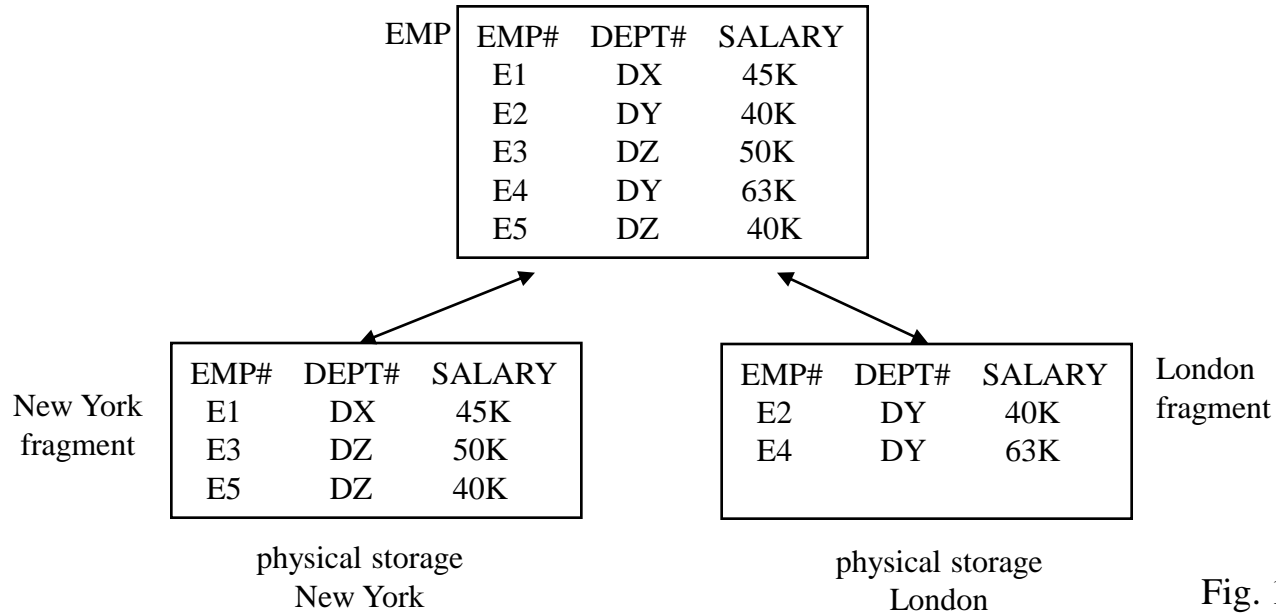


Fig. 12.2: An example of fragmentation.

The Twelve Objectives (cont.)

- Data Fragmentation
 - A fragmentation can be any subrelation derivable via restriction and projection (with primary key).
 - Advantage: data can stored at the location where it is most frequently used.
- **Fragmentation independence**
 - user should be able to behave as if the relations were not fragmented at all.
 - one reasons why relational technology is suitable for DBMS.
 - user should be presented with a view of data.
=> system must support updates against join and union views.
 - Advantages
 - (1) simplify user program and activity.
 - (2) allow data to be re-fragmented at any time.

The Twelve Objectives (cont.)

★ 6. Replication Independence (Replication Transparency)

- Data Replication

USER PERCEPTION

EMP	EMP#	DEPT#	SALARY
	E1	DX	45K
	E2	DY	40K
	E3	DZ	50K
	E4	DY	63K
	E5	DZ	40K

New York fragment

EMP#	DEPT#	SALARY
E1	DX	45K
E3	DZ	50K
E5	DZ	40K

London fragment

EMP#	DEPT#	SALARY
E2	DY	40K
E4	DY	63K

replica of London fragment

EMP#	DEPT#	SALARY
E2	DY	40K
E4	DY	63K

copy

replica of New York fragment

EMP#	DEPT#	SALARY
E1	DX	45K
E3	DZ	50K
E5	DZ	40K

physical storage
New York

physical storage
London

Fig. 12.3: An example of replication.

The Twelve Objectives (cont.)

- **Data Replication**
 - A given fragment of relation can be represented at the physical level by many distinct copies of the same object at many distinct sites.
 - Unit of replication: fragment (may not a complete relation)
 - Advantage: better performance and availability
 - Disadvantage: update propagation problem.
- **Replication Independence**
 - User should be able to behave as if the data is not replicated at all.
 - Advantages
 - (1) simplify user programs and activities.
 - (2) allow replicas to be created and destroyed dynamically.

The Twelve Objectives (cont.)

7. Distributed Query Processing

- message transfer cost
- optimization

8. Distributed Transaction Management

- concurrency control
- recovery control

9. Hardware Independence: IBM, DEC, HP, PC, ...

10. Operating System Independence: VMS, UNIX, ...

11. Network Independence: BITNET, INTERNET, ARPANET, ...

12. DBMS Independence: Relational, hierarchical, network, ...

- distributed system may be heterogeneous.

12.3 Problems of Distributed Database Systems

Basic Point: Network are slow !

Basic point: network are slow !



Overriding Objective : *minimize the number and volume of messages.*



Give rise to the following problem

- Query Processing
- Update Propagation
- Concurrency
- Recovery
- Catalog Management

⋮

Query Processing: Example (cont.)

- Query: "Select S# for London suppliers of Red Parts"

SELECT	S.S#	site A	site B
FROM	S, P, SP	S, SP	P
WHERE	S.CITY = "London"		
	AND S.S# = SP.S#	S \bowtie SP	
	AND SP.P# = P.P#		
	AND P.COLOR = 'Red'		

- Estimates
 - # of Red Parts = 10
 - # of Shipments by London Supplier = 100,000
- Communication Assumption :
 - Data Rate = 10,000 bits per second
 - Access Delay = 1 second
- $T[i]$ = total communication time for strategy i
 - = total access delay + total data volume / data rate
 - = (# of messages * 1 sec) + (total # of bits / 10,000) sec.

Query Processing: Example (cont.)

- Strategy 1

1. Join S and SP at site A

site A
S, SP

site B
P

2. Select tuples from ($S \bowtie SP$) for which city is 'London'
(100,000 tuples)

3. For each of those tuple, check site B to see if the part is red. (2 messages: 1 query, 1 response)

$$T[1] = (100,000 * 2) * 1 = \mathbf{2.3 \text{ days}}$$

- Strategy 2

Move relations S and SP to site B and process the query at B.

$$T[2] = 2 + (10,000 + 1,000,000) * 100 / 10,000 = \mathbf{28 \text{ hours}}$$

- Strategy 3

Move relation P to site A and process the query at A

$$T[3] = 1 + (100,000 * 100) / 10,000 = \mathbf{16.7 \text{ min}}$$

Query Processing: Example (cont.)

- Strategy 4

1. Select tuples from P where color is red. (10 tuples)
2. Check site A to see if there exists a shipment relating the part to a London Supplier. (2*10 messages)

$$T[4] = 2 * 10 * 1 = 20 \text{ sec}$$

site A
S, SP

site B
P

- Strategy 5

1. Select tuples from P where color is red (10 tuples)
2. Move the result to site A and complete the processing at A.

$$T[5] = 1 + (10 * 100) / 100,000 = 1.01 \text{ sec}$$

- *Note:* Each of the five strategies represents a plausible solution, but the variation in communication time is enormous.

Query Processing: Semijoin

- Semijoin: (used in SDD - 1) Ref. p.529 [18.15]

$A \bowtie B$
 \bowtie
 \bowtie

p.626 [21.26]

<e.g.>

- Database :

S: 1,000 tuples, at site A

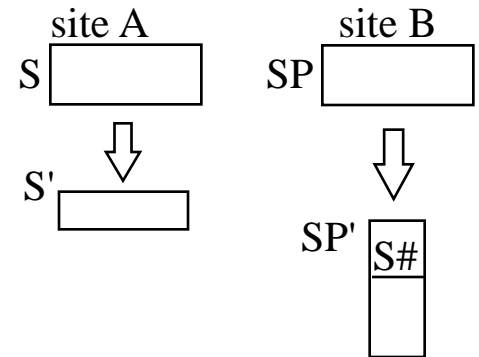
SP: 2,000 tuples, at site B

of tuples in S where S.S#=SP.S#: 100,

length of a S tuple: 100 bit

length of a SP tuple: 100 bit

length of the S# field: 10 bit



- Regular Join:

<1> Ship S to site B (1000 * 100 bits)

<2> Join S and SP at site B

communication time = 1 + 1000*100/10000 = 11 sec

Query Processing: Semijoin (cont.)

- Semijoin

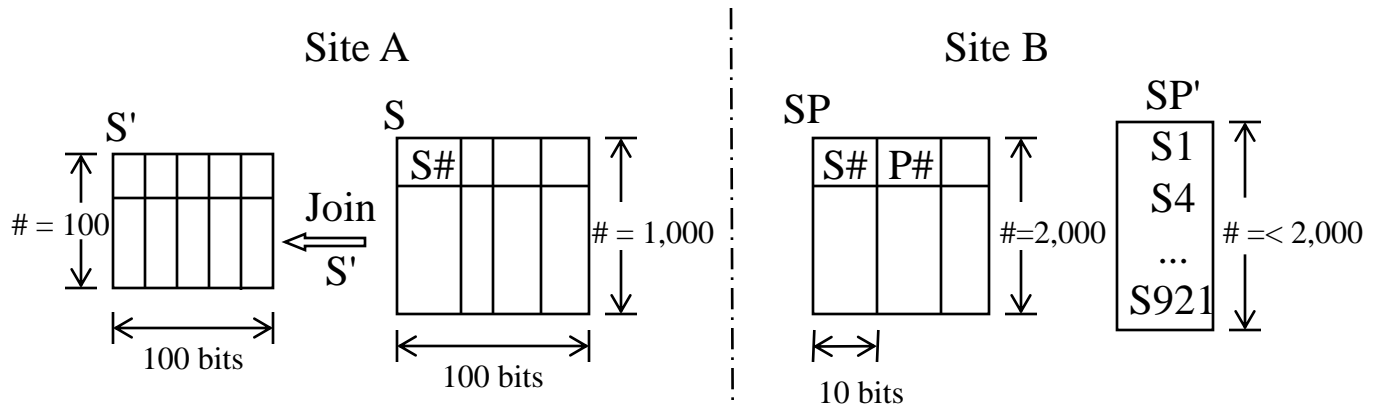
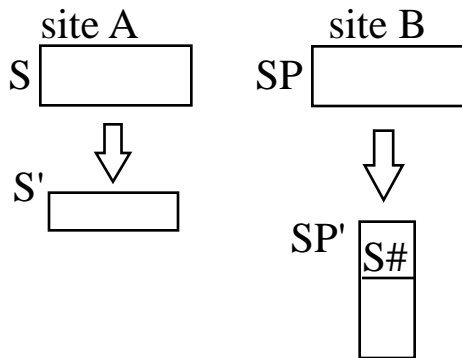
<1> site B: step 1. Project SP on S# (get SP')
 step 2. ship to site A

<2> site A: step 3. Join the projection of SP' on S# with S
 step 4. The result S', ship to site B

<3> site B: step 5. Join S' with SP

$$\text{communication time} = 1 + 10 * 2000 / 10000 + 1 + 100 * 100 / 10000$$

$$= 1 + 2 + 1 + 1 = \mathbf{5 \text{ sec}}$$

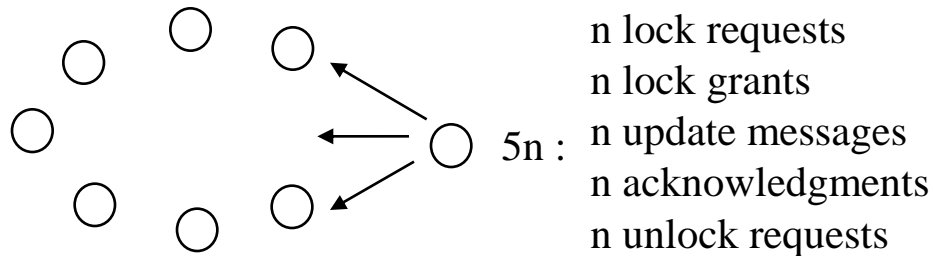


Update Propagation

- Basic problem with data replication
 - An update to any given logical data object must be propagated to all stored copies of that object.
 - some sites may be unavailable (because of site or network failure) at the time of the update
=> Data is less available !
- A possible Solution: Primary Copy (used in distributed INGRES)
 - one copy of each object is designated as the *primary copy*.
 - primary copies of different objects will generally be at different sites.
 - Update Operation
 1. Complete as soon as the primary copy has been updated.
 2. Control is returned and the transaction can continue execute.
 3. The site holding the primary copy broadcasts the update to all other sites.
 - Further Problem: violation of the local autonomy objective.

Concurrency

- Most distributed systems are based on locking .
- Requests to *test*, *set* and *release* locks are *messages* → overhead.
<e.g.>: Consider a transaction **T** that needs to update an object for which there exists replicas at n remote sites.



- Several orders of magnitude greater than in a centralized system.
- Solution
 - adopt the primary copy strategy.
 - the site holding the primary copy of X handles all locking operations involving X.
 - 1 lock request, 1 lock grant, n updates, n ack, and 1 unlock request ($2n+3 \leq 5n$).
 - Problem : loss of local autonomy.

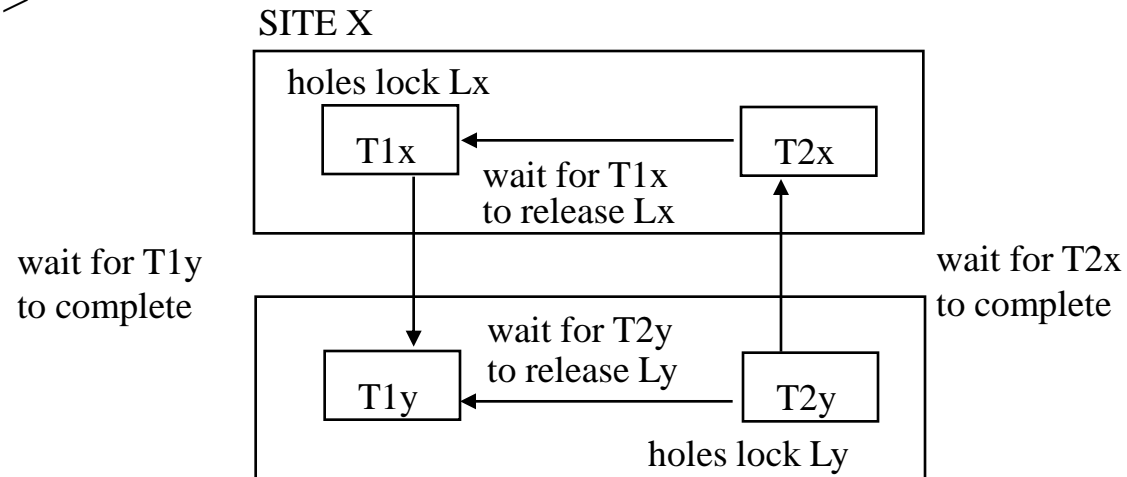
Concurrency (cont.)

- **Global Deadlock Problem**

- Neither site can detect it using only information that is internal to that site.

i.e. no cycles in two local wait-for-graph, but a cycle in the global.

<e.g.>



- Global deadlock detection needs further communication overhead.

Recovery

■ Two-Phase Commit

- important whenever a given transaction can interact with multiple, independent 'resource manager' (site).
- the transaction issues a single system-wide COMMIT (or ROLLBACK), which is handled by a new system component: Coordinator.
- The coordinator goes through the following two-phase process :
 - <1> Coordinator requests all participant to decide whether commit (reply OK) or rollback.
 - <2> If all participants reply OK, the coordinator broadcast 'COMMIT', otherwise broadcast 'ROLLBACK' , and all participants must obey.

■ Points arising

- the coordinator function must be performed by different sites for different transactions => No reliance on a central site. (P. 612)
- the participants must do what is told by the coordinator => loss of local autonomy.
- Coordinator must communicate with participants => more messages and more overhead.
- No protocol guarantee that all participants will commit or rollback in unison.

Catalog Management

- Contents of catalog:
 - not only data regarding relations, indexes, users, etc,
 - but also all the necessary control information for independence.
- Where and How ?
 1. **Centralized:** violate no reliance on a central site.
 2. **Fully replicated:** loss of local autonomy.
 3. **Partitioned:** nonlocal operations are very expensive.
 4. **Combination** of (1) and (3): violate no reliance on a central site.

12.4 Gateways

Gateways

- DBMS independence
- Suppose INGRES provides a "gateway", that "making ORACLE look like INGRES"

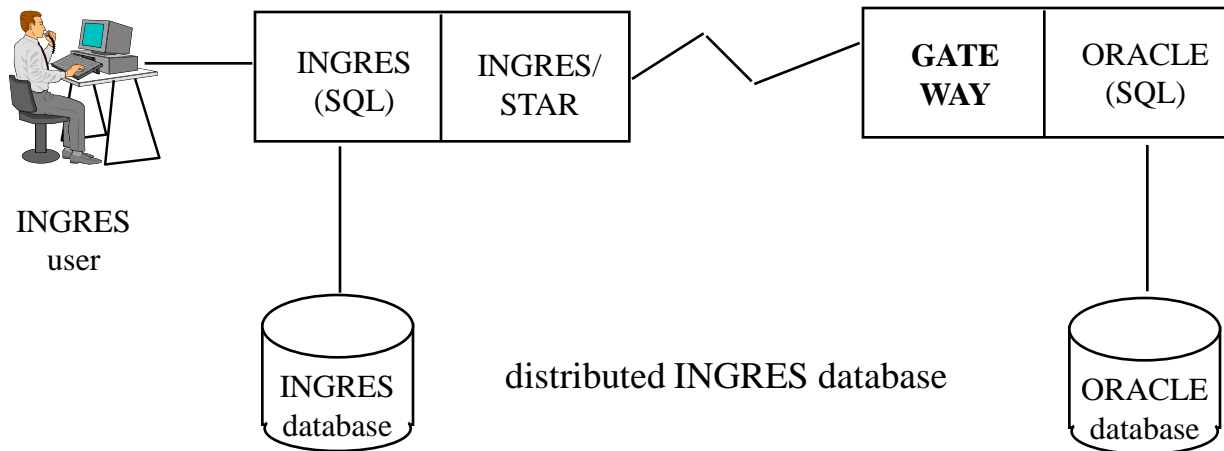


Fig. 12.4: A hypothetical INGRES-provided gateway to ORACLE

12.5 Client/Server Systems

Client/Server Architecture

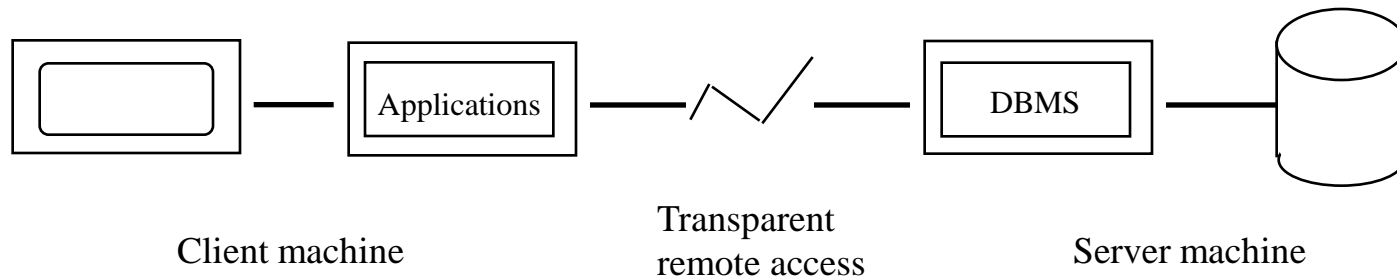


Fig. 12.5: A client/server system

- some sites are client, and others are server sites
- a great deal of commercial products
- little in "true" general-purpose distributed system (but long-term trend might be important)
- client: application or front-end
- server: DBMS or backend
- Several variations

Client/Server Systems

■ Client/Server Standards

- SQL/92 standard
- ISO (International Organization for Standardization)
- Remote Data Access standard: RDA.
- Distributed Relational Database Architecture standard: DRDA.

■ Client/Server Application Programming

- Stored procedure: a precompiled program that is stored at the server site.
- invoked from the client by a remote procedure call (RPC).
- one stored procedure can be shared by many clients.
- optimization can be done at the precompiled time instead of at run time.
- provide better security.
- but, no standards in this area.

ADVANCED DATABASE SYSTEMS

國立東華大學 資訊管理系
National Dong Hwa University
Department of Information Management

楊維邦 教授
Prof. Weipang Yang

Contents: Part I

- ❑ **Unit 0** Introduction to Database Systems
- ❑ **Unit 1** Introduction to DBMS
- ❑ **Unit 2** DB2 and SQL
- ❑ **Unit 3** The Relational Model
- ❑ **Unit 4** The Hierarchical Model
- ❑ **Unit 5** The Network Model
- ❑ **Unit 6** Access Methods

- ❑ **References:**

1. C. J. Date, An Introduction to Database Systems, 8th edition, 2004.
2. J. D. Ullman, Principles of Database and Knowledge-Base, Vol.I, 1988.
3. Cited papers

Contents: Part II

- ❑ **Unit 7 Logical Database Design**
- ❑ **Unit 8 Database Recovery**
- ❑ **Unit 9 Concurrency Control**
- ❑ **Unit 10 Security and Integrity**
- ❑ **Unit 11 Query Optimization**
- ❑ **Unit 12 Distributed Database**

- ❑ **References:**

1. C. J. Date, *An Introduction to Database Systems*, 8th edition, 2004.
2. J. D. Ullman, *Principles of Database and Knowledge-Base*, Vol. I, 1988.
3. Cited papers

Contents: Part III

- ❑ **Unit 13** **Object-oriented Database**
- ❑ **Unit 14** **Logic-Based Database**
- ❑ **Unit 15** **Image Database**
- ❑ **Unit 16** **Multimedia Database**
- ❑ **Unit 17** **Multidatabase**
- ❑ **Unit 18** **Real-time Database**
- ❑ **Unit 19** **Parallel Database**
- ❑ **Unit 20** **Temporal Database**
- ❑ **Unit 21** **Active Database**
- ❑ **Unit 22** **Bioinformatic Database**
- ❑ **Unit 23** **....**

Study and Research on Databases

Level 5: Doing Research

Level 4: Survey Papers: Special Topics (Unit 13 -)

Level 3: DBMS: Advanced Topics

(Unit 7 – 12)

Date, Vol. 1, 2

Ullman

Level 2: DBMS: Fundamentals

(Unit 1 – 6)

Date, Vol. 1

Using DB2,

Project: Design a “mini” DBMS

Level 1: Using DBMS

Where to find Database papers?

□ Journals:

- ACM Trans. On Database Systems (TODS)
- IEEE Trans. On Software Engineering (TOSE)
- IEEE Trans. On Knowledge and Data Engineering
- ACM Trans. On Office Information Systems (TOOIS)
- **Journal of Information Science and Engineering**
- ...

□ Conferences:

- VLDB (Very Large Data Base)
- ACM SIGMOD
- PODS (Principal of Database Systems)
- IEEE Data Engineering
- **NCS, ICS (National/International Computer Symposium)**
- ...