

DISTRIBUTED DATABASES

CS561-SPRING 2012
WPI, MOHAMED ELTABAKH

1

RECAP: PARALLEL DATABASES

- **Three possible architectures**
 - Shared-memory
 - Shared-disk
 - Shared-nothing (the most common one)
- **Parallel algorithms**
 - **Intra-operator**
 - Scans, projections, joins, sorting, set operators, etc.
 - **Inter-operator**
 - Distributing different operators in a complex query to different nodes
- **Partitioning and data layout is important and affect the performance**
 - Range-based, hash-based, round robin
- **Optimization of parallel algorithms is a challenge**

DISTRIBUTED DATABASE

DEFINITIONS

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) = DB + Communication

DISTRIBUTED DATABASES

MAIN CONCEPTS

- **Data are stored at several locations**
 - Each managed by a DBMS that can run autonomously
- **Ideally, location of data is unknown to client**
 - ***Distributed Data Independence***
- **Distributed Transactions**
 - Clients can write Transactions regardless of where the affected data are located
 - ***Big question:*** How to ensure the ACID properties Distributed Transactions???

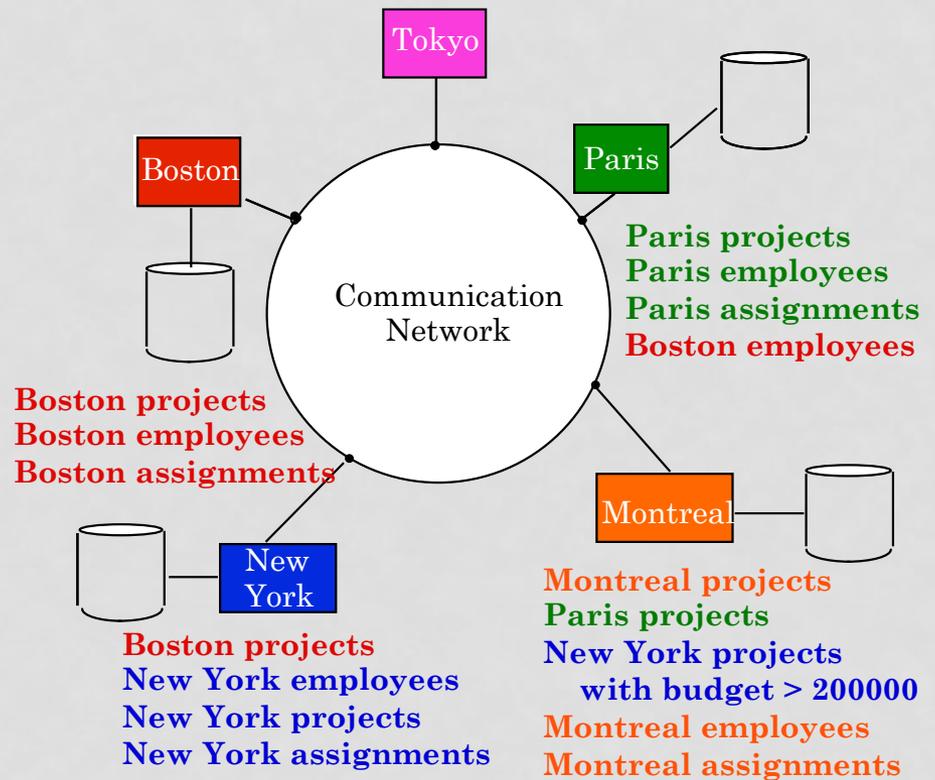
DISTRIBUTED DBMS PROMISES

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

TRANSPARENCY & DATA INDEPENDENCE

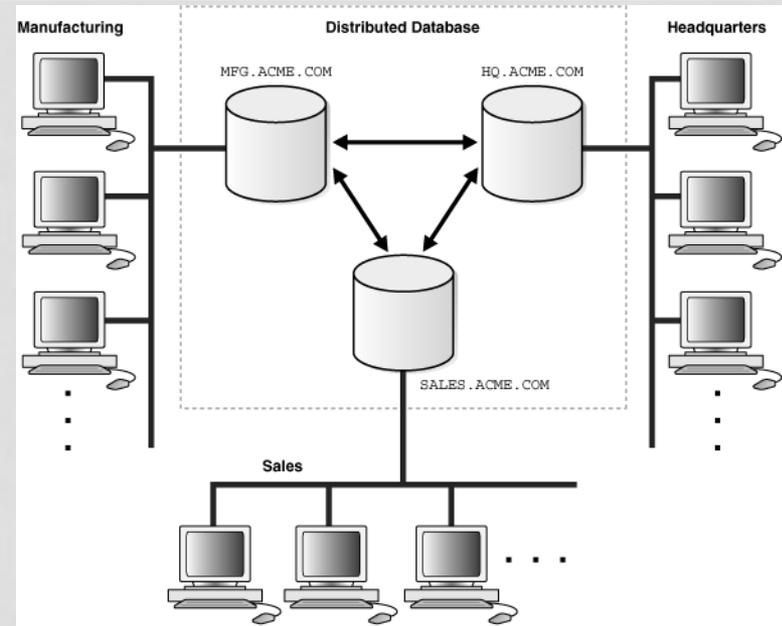
- Data distributed (with some replication)
- Transparently ask query:

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

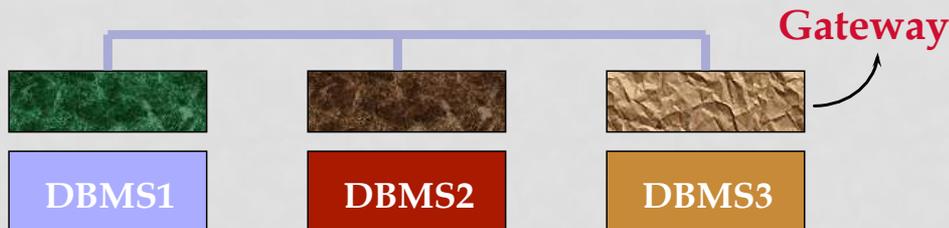


TYPES OF DISTRIBUTED DATABASES

- **Homogeneous**
 - Every site runs the same type of DBMS
- **Heterogeneous:**
 - Different sites run different DBMS (maybe even RDBMS and ODBMS)



Homogeneous DBs can communicate directly with each other



Heterogeneous DBs communicate through gateway interfaces

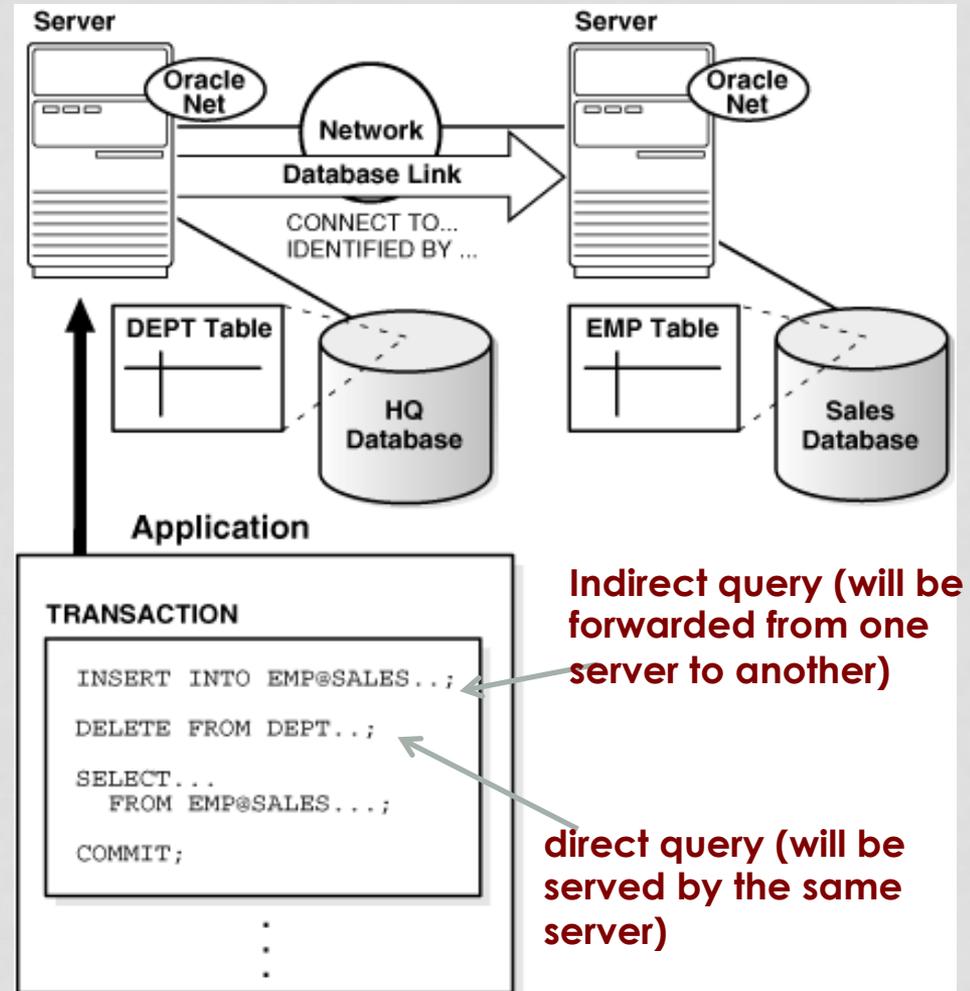
DISTRIBUTED DATABASE ARCHITECTURE

- **Client-Server**

- Client connects directly to specific server(s) and access only their data
- Direct queries only

- **Collaborative Servers**

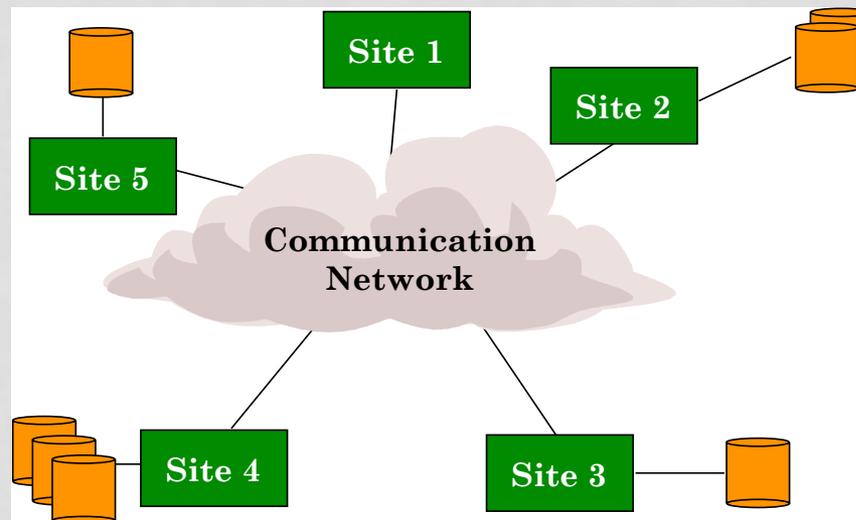
- Servers can serve queries or be clients and query other servers
- Support indirect queries



DISTRIBUTED DATABASE ARCHITECTURE (CONT'D)

- **Peer-to-Peer Architecture**

- Scalability and flexibility in growing and shrinking
- All nodes have the same role and functionality
- Harder to manage because all machines are **autonomous** and **loosely coupled**



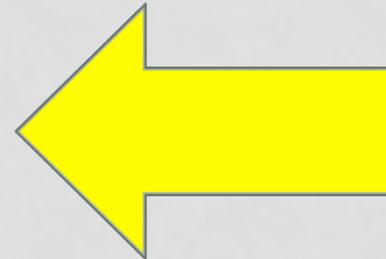
MAIN ISSUES

- **Data Layout Issues**
 - Data partitioning and fragmentation
 - Data replication
- **Query Processing and Distributed Transactions**
 - Distributed join
 - Transaction atomicity using two-phase commit
 - Transaction serializability using distributed locking

MAIN ISSUES

- **Data Layout Issues**

- Data partitioning and fragmentation
- Data replication



- **Query Processing and Distributed Transactions**

- Distributed join
- Transaction atomicity using two-phase commit
- Transaction serializability using distributed locking

FRAGMENTATION

- **How to divide the data? Can't we just distribute relations?**
- **What is a reasonable unit of distribution?**
 - **relation**
 - views are subsets of relations
 - extra communication
 - Less parallelism
 - **fragments of relations (sub-relations)**
 - concurrent execution of a number of transactions that access different portions of a relation
 - views that cannot be defined on a single fragment will require extra processing
 - semantic data control (especially integrity enforcement) more difficult

FRAGMENTATION ALTERNATIVES – HORIZONTAL

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

PROJ₁

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

Stored in London

PROJ₂

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

Stored in Boston

FRAGMENTATION ALTERNATIVES – VERTICAL

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

Horizontal partitioning is
more common

PROJ₁

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

Stored in London

PROJ₂

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

Stored in Boston

CORRECTNESS OF FRAGMENTATION

- **Completeness**

- Decomposition of relation R into fragments R_1, R_2, \dots, R_n is complete if and only if each data item in R can also be found in some R_i

- **Reconstruction (Lossless)**

- 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 (Non-overlapping)**

- 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 ($k \neq j$).

REPLICATION 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

DATA REPLICATION

- **Pros:**

- Improves availability
- Disconnected (mobile) operation
- Distributes load
- Reads are cheaper

- **Cons:**

- N times more updates
- N times more storage

- **Synchronous vs. asynchronous**

- **Synchronous:** all replica are up-to-date
- **Asynchronous:** cheaper but delay in synchronization

Catalog Management

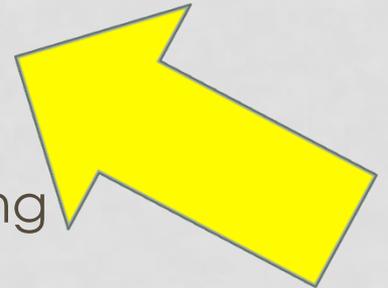
- Catalog is needed to keep track of the location of each fragment & replica
- Catalog itself can be centralized or distributed

COMPARISON OF REPLICATION ALTERNATIVES

	Full-replication	Partial-replication	Partitioning
QUERY PROCESSING	Easy	← Same Difficulty →	
DIRECTORY MANAGEMENT	Easy or Non-existent	← Same Difficulty →	
CONCURRENCY CONTROL	Moderate	Difficult	Easy
RELIABILITY	Very high	High	Low
REALITY	Possible application	Realistic	Possible application

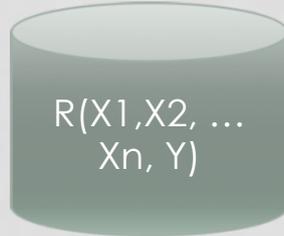
MAIN ISSUES

- **Data Layout Issues**
 - Data partitioning and fragmentation
 - Data replication
- **Query Processing and Distributed Transactions**
 - Distributed join
 - Transaction atomicity using two-phase commit
 - Transaction serializability using distributed locking



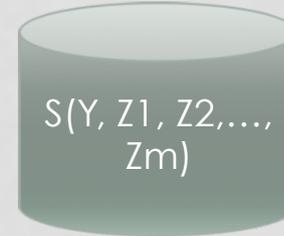
DISTRIBUTED JOIN $R(X,Y) \bowtie S(Y,Z)$

Stored in London



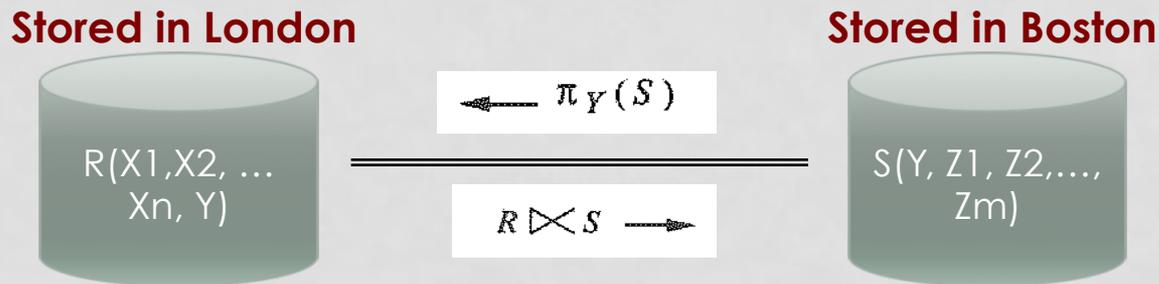
Join based on
 $R.Y = S.Y$

Stored in Boston



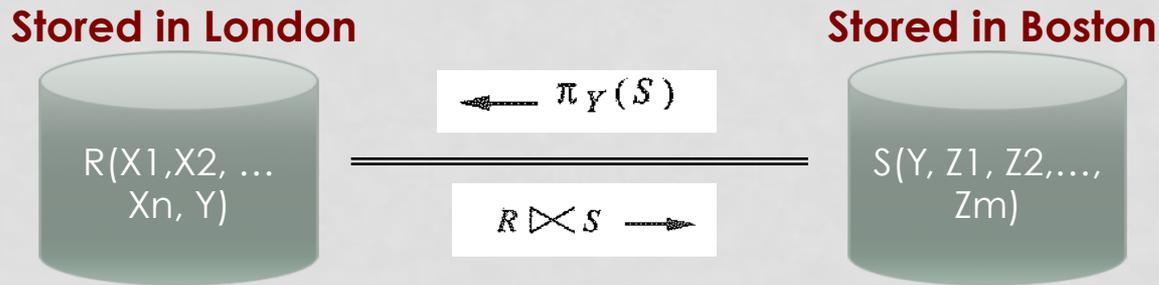
- **Option 1:** Send R to S's location and join their
- **Option 2:** Send S to R's location and join their
- *Communication cost is expensive, too much data to send*
- **Is there a better option ???**
 - **Semi Join**
 - **Bloom Join**

SEMI-JOIN



- Send only S.Y column to R's location
- Do the join based on Y columns in R's location (**Semi Join**)
- Send the records of R that will join (without duplicates) to S's location
- Perform the final join in S's location

IS SEMI-JOIN EFFECTIVE



Depends on many factors:

- If the size of Y attribute is small compared to the remaining attributes in R and S
- If the join selectivity is high $\rightarrow R \bowtie S$ is small
- If there are many duplicates that can be eliminated

BLOOM JOIN

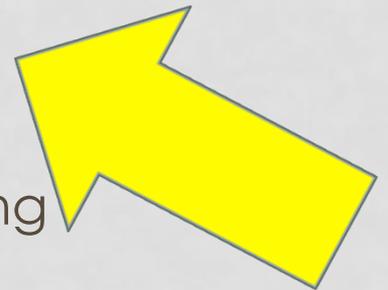
- Build a bit vector of size K in R's location (all 0's)



- **For every record in R, use a hash function(s) based on Y value (return from 1 to K)**
 - Each function hashes Y to a bit in the bit vector. Set this bit to 1
- Send the bit vector to S's location
- **S will use the same hash function(s) to hash its Y values**
 - If the hashing matched with 1's in all its hashing positions, then this Y is candidate for Join
 - Otherwise, not candidate for join
 - Send S's records having candidate Y's to R's location for join

MAIN ISSUES

- **Data Layout Issues**
 - Data partitioning and fragmentation
 - Data replication
- **Query Processing and Distributed Transactions**
 - Distributed join
 - **Transaction atomicity using two-phase commit**
 - Transaction serializability using distributed locking



TRANSACTIONS

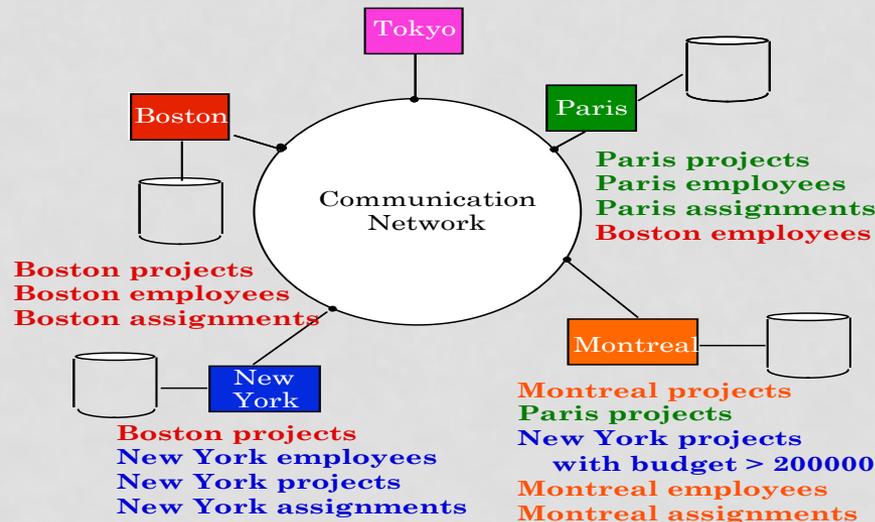
- A Transaction is an atomic sequence of actions in the Database (reads and writes)
- Each Transaction has to be executed *completely*, and must leave the Database in a consistent state
- If the Transaction fails or aborts midway, then the Database is “rolled back” to its initial consistent state (before the Transaction began)



ACID Properties of Transactions

ATOMICITY IN DISTRIBUTED DBS

- **One transaction T may touch many sites**
 - T has several components T_1, T_2, \dots, T_m
 - Each T_k is running (reading and writing) at site k
 - **How to make T is atomic ????**
 - Either T_1, T_2, \dots, T_m complete or None of them is executed
- **Two-Phase Commit techniques is used**

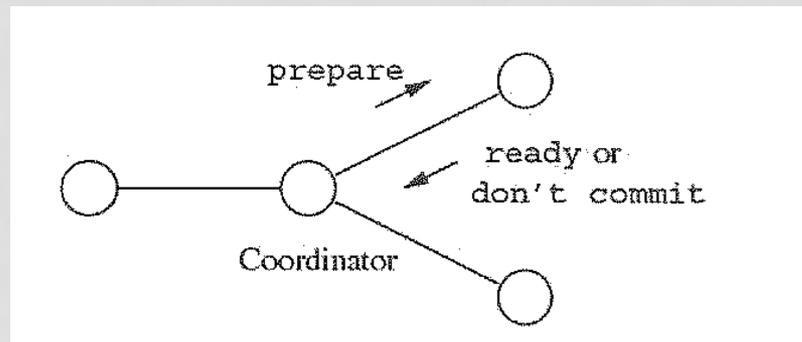


TWO-PHASE COMMIT

- **Phase 1**

- Site that initiates T is the **coordinator**
- When coordinator wants to commit (complete T), it sends a “*prepare T*” msg to all participant sites
- Every other site receiving “*prepare T*”, either sends “*ready T*” or “*don't commit T*”
 - A site can wait for a while until it reaches a decision (Coordinator will wait reasonable time to hear from the others)

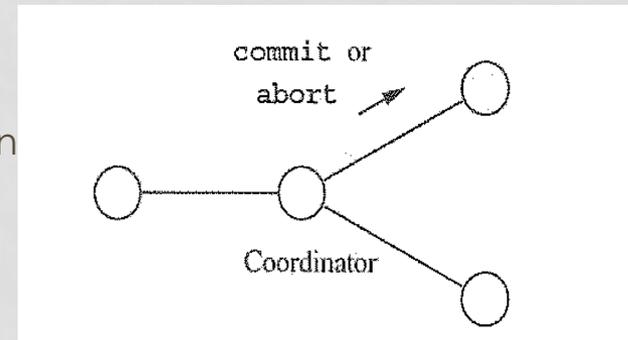
- **These msgs are written to local logs**



TWO-PHASE COMMIT (CONT'D)

- Phase 2

- **IF coordinator received all “ready T”**
 - Remember no one committed yet
 - Coordinator sends “*commit T*” to all participant sites
 - Every site receiving “*commit T*” commits its transaction
- **IF coordinator received any “don’t commit T”**
 - Coordinator sends “*abort T*” to all participant sites
 - Every site receiving “*abort T*” commits its transaction



- These msgs are written to local logs

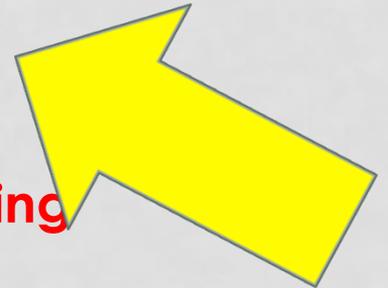
Example 1: What if one sites in Phase 1 replied “don’t commit T”, and then crashed???

Example 2: What if all sites in Phase 1 replied “ready T”, then one site crashed???

- Straightforward if no failures happen
- In case of failure logs are used to ensure **ALL** are done or **NONE**

MAIN ISSUES

- **Data Layout Issues**
 - Data partitioning and fragmentation
 - Data replication
- **Query Processing and Distributed Transactions**
 - Distributed join
 - Transaction atomicity using two-phase commit
 - **Transaction serializability using distributed locking**



DATABASE LOCKING

- Locking mechanisms are used to prevent concurrent transactions from updating the same data at the same time
- Reading(x) → shared lock on x
- Writing(x) → exclusive lock on x
- **More types of locks exist for efficiency**

What you request

What you have

	Shared lock	Exclusive lock
Shared lock	Yes	No
Exclusive lock	No	No

In Distributed DBs:

- x may be replicated in multiple sites (not one place)
- The transactions reading or writing x may be running at different sites

DISTRIBUTED LOCKING

- **Centralized approach**
 - One dedicated site managing all locks
 - Cons: bottleneck, not scalable, single point of failure
- **Primary-Copy approach**
 - Every item in the database, say x , has a primary site, say P_x
 - Any transaction running any where, will ask P_x for lock on x
- **Fully Distributed approach**
 - To read, lock any copy of x
 - To write, lock all copies of x
 - Variations exists to balance the cots of read and write op.

Deadlocks are very possible. How to resolve them???

Using timeout: After waiting for a while for a lock, abort and start again

SUMMARY OF DISTRIBUTED DBS

- ***Promises of DDBMSs***

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

- ***Classification of DDBMS***

- Homogeneous vs. Heterogeneous
- Client-Server vs. Collaborative Servers vs. Peer-to-Peer

SUMMARY OF DISTRIBUTED DBS (CONT'D)

- **Data Layout Issues**
 - Data partitioning and fragmentation
 - Data replication
- **Query Processing and Distributed Transactions**
 - Distributed join
 - Transaction atomicity using two-phase commit
 - Transaction serializability using distributed locking