

Distributed Databases

COSC430—Advanced Databases David Eyers

Learning objectives

You should be able to

- define the key concepts in distributed database
- distinguish between different types of distributed database
- understand the architectures of distributed database
- explain how to perform data fragmentation, allocation, and replication
- explain how to do semi-join in distributed database
- understand technologies for distributed transactions
- You will get a taste of carrying out scientific research
 - optimisation on data allocation and replication
 - the study of the research paper on Google Spanner

COSC430 Lecture 5, 2020

ributed database ypes of distributed database of distributed database ragmentation, allocation, and

distributed database istributed transactions ng out scientific research n and replication er on Google Spanner

2

Definitions

- Distributed database
 - distributed over a **network** of computers
- Distributed DBMS
 - transparent to the users

COSC430 Lecture 5, 2020

a collection of multiple, logically interrelated databases

 a software system that manages a distributed database and provides an access mechanism while making the distribution

Distributed database system=Database+Communication





Key concepts for distributed databases

- Data stored at several locations
 - Fragmentation
 - Replication EMPLOYEES
- Transparency: hide implementation for users/developers
- PROJECTS WORKS_ON

PROJECTS

- Location
- Fragmentation
- Replication
- Design
- Execution

COSC430 Lecture 5, 2020—reproduction of Figure 25.1 from E&N





Types of distributed databases

- Homogeneous
 - All sites run the same DBMS

 Heterogeneous Different sites can run different DBMSs





Distributed database architecture

- Client-Server
 - a two-level architecture based on the client-server concept (database centric)
 - A client can directly or indirectly connect to a server



Distributed Database Architecture

- Peer-to-Peer
 - All nodes have the same role and functionality
 - High scalability and flexibility
 - Servers are autonomous and loosely coupled
 - Management layer determine whether such infrastructure is easier or harder to manage







Main database functionality to consider

- Data layout
 - Data fragmentation
 - Data allocation and replication
- Query processing and optimisation
 - Data transfer cost
 - Semi-join
 - Query and update decomposition
- Distributed transactions

 - Transaction atomicity using two-phase commit Transaction serialisability using distributed locking



Data fragmentation

- which may be assigned for storage at various sites
- What is a reasonable unit of data distribution?
 - Relation
 - Can increase communications
 - Less parallelism
 - Sub-relations

 - Better parallelism
 - Difficult to enforce integrity constraints

COSC430 Lecture 5, 2020

Break up the database into logical units, called fragments,

Extra processing for views that cannot be defined on a single fragment



Data Fragmentation

Horizontal fragmentation

- Grouping rows to create subsets of tuples
- Fach subset has a certain

					PNUMBER	PLOCATION	DNUM	
logical meaning				Reorganisation		Houston	1	
PNUMBER	PLOCATION	DNUM		PNAME	PNUMBE	R PLOCATION	DI	
20	Houston	1		Computerisation	1	0 Stafford		
10	Stafford	4	\longrightarrow	Newbenefits	3	0 Stafford		
30	Stafford	4	•	NonProject	4	0		
40		4						
1	Bellaire	5						
3	Houston	5		PNAME	PNUMBER	PLOCATION	DNUM	
2	Sugarland	5						
				ProductX	1	Bellaire	5	
				ProductY	2	Sugarland	5	
				ProductZ	3	Houston	5	
	20 PNUMBER 20 10 30 40 1 3 2	PNUMBER PLOCATION 20 Houston 10 Stafford 30 Stafford 40 1 Bellaire 3 Houston 2 Sugarland	PNUMBER PLOCATION DNUM 20 Houston 1 10 Stafford 4 30 Stafford 4 40 4 1 Bellaire 5 3 Houston 5 2 Sugarland 5	PNUMBER PLOCATION DNUM 20 Houston 1 10 Stafford 4 30 Stafford 4 40 1 Bellaire 5 3 Houston 5 2 Sugarland 5	PNUMBER PLOCATION DNUM 20 Houston 1 10 Stafford 4 30 Stafford 4 40 4 1 Bellaire 5 3 Houston 5 2 Sugarland 5	PNUMBER PLOCATION DNUM 20 Houston 1 20 Houston 1 30 Stafford 4 40 4 1 Bellaire 5 3 Houston 5 2 Sugarland 5	PNAME PNUMBER PLOCATION 20 Houston 1 20 Houston 1 10 Stafford 4 30 Stafford 4 40 4 1 Bellaire 5 3 Houston 5 2 Sugarland 5	







Data Fragmentation

 Vertical fragmentation—divide relation by columns that the full relation can be reconstructed

FNAME	М	LNAME	SSN	BDATE	ADDRESS	S	SALARY	SUPERSSN	DNO
 John	 R	Smith	123456789	09-14N-65	731 Fondren, Houston, TX		30000	333445555	5
Franklin	Т	Wong	333445555	08-DEC-55	638 Voss, Houston, TX	M	40000	888665555	5
Alicia	J	Zelaya	999887777	19–JUL–68	3321 Castle, Spring,TX	F	25000	987654321	4
Jennifer	S	Wallace	987654321	20-JUN-41	291 Berry, Bellaire, TX	F	43000	888665555	4
Ramesh	K	Narayan	666884444	15-SEP-62	975 Fire Oak, Humble, TX	M	38000	333445555	5
Joyce	Α	English	453453453	31-JUL-72	5631 Rice, Houston, TX	F	25000	333445555	5
Ahmad	v	Jabbar	987987987	29-MAR-69	980 Dallas, Houston, TX	M	25000	987654321	4
James	E	Borg	888665555	10-NOV-37	450 Stone, Houston, TX	М	55000		1

FNAME	М	LNAME	SSN	BDATE	ADDRESS
	-				
John	В	Smith	123456789	09-JAN-65	731 Fondren, Houston, T
Franklin	т	Wong	333445555	08-DEC-55	638 Voss, Houston, TX
Alicia	J	Zelaya	999887777	19–JUL–68	3321 Castle, Spring,TX
Jennifer	S	Wallace	987654321	20-JUN-41	291 Berry, Bellaire, TX
Ramesh	κ	Narayan	666884444	15-SEP-62	975 Fire Oak, Humble, T
Joyce	Α	English	453453453	31-JUL-72	5631 Rice, Houston, TX
Ahmad	۷	Jabbar	987987987	29-MAR-69	980 Dallas, Houston, TX
James	Е	Borg	888665555	10-N0V-37	450 Stone, Houston, TX

COSC430 Lecture 5, 2020

Each fragment has the primary key or some candidate key so



SSN	SALARY	SUPERSSN	DNO
123456789	30000	333445555	5
333445555	40000	888665555	5
999887777	25000	987654321	4
987654321	43000	888665555	4
666884444	38000	333445555	5
453453453	25000	333445555	5
987987987	25000	987654321	4
888665555	55000		1



Data Fragmentation

- Mixed (Hybrid) fragmentation
 - Intermix the two types of fragmentation
- Correctness of fragmentation
 - Completeness: each data item can be found in one fragment
 - Reconstruction (lossless): the full relation can be reconstructed from all fragments
 - Disjointness (non-overlapping): each data item except the key should not be in more than one fragment





Data Replication

- Non-replicated
 - Each fragment resides at only one site
- Replicated
 - Fully replicated: each fragment at each site
 - Partially replicated: each fragments at some sites
- Pros & cons
 - Improve availability, distribute load, cheaper reads
 - Complexity on update and storage
- Rule of thumb
 - otherwise replication may cause problems

COSC430 Lecture 5, 2020

• If (read-only queries)/(update queries) \geq 1, replication is advantageous,



Comparison of replication alternatives

	Full replication	Partial replication	Partitioning		
Query processing	Easy	Some difficulty			
Directory management	Easy or non- existent	Some difficulty			
Concurrency control	Moderate	Difficult	Easy		
Reliability	Very high	High	Low		



Optimisation

- Best fragmentation, replication and allocation?
 - database system (Corcoran, SAC'94)
 - Set of m sites S, each has capacity ci
 - Set F of n fragments, each fragment j has size si
 - Site requirement matrix R:
 - r_{i,i} is requirement by site i for fragment j

A generic algorithm for fragment allocation in a distributed

 $S = \{c_1, c_2, c_3, \dots, c_i, \dots, c_m\}$

 $F = \{s_1, s_2, s_3, \dots, s_i, \dots, s_n\}$

R

ļ	$[r_{1,1}]$	<i>r</i> _{1,2}	• • •	$r_{1,n}$
	<i>r</i> _{2,1}	$r_{2,2}$	• • •	$r_{2,n}$
_	•	• •	•••	• •
	$r_{m,1}$	$r_{m,2}$	• • •	$r_{m,n}$



Formalising optimisation of allocation

- Transmission cost matrix T
- Fragment placement vector $P = \{p_1, p_2, p_3, ..., p_j, ..., p_n\}$ p_i=i indicates fragment j is at site i
- n • Objective: minimise total transmission cost: $\sum \sum r_{i,j} t_{i,p_i}$ • ... but subject to: $\forall i, 1 \leq i \leq m$, $\sum_{i,j}^{n} r_{i,j} s_j \leq c_i$ i=1 j=11 = 1





Distributed query processing

- - Option 1: Send both R and S to Site 1 for join
- Which option is best? Is it optimal?



COSC430 Lecture 5, 2020

Optimisation goal: reduce the amount of data transfer Option 2: Send R to Site 3 to join, send join results back to site 1 Option 3: Send S to Site 2 to join, send join results back to site 1

Site 1

$$R(X_1, X_2, ..., X_n, Y)$$
 Site 2
 $S(Y, Z_1, Z_2, ..., Z_n)$ Site 3



Semi-joins can improve efficiency of queries

- Reduce number of tuples before transfer to other site Site 3 sends only S.Y column to Site 2

 - Site 2 does the join based on R's Y column; sends the records of R that will join (without duplicates) back to Site 3
 - Site 3 performs the final join









Semi-join example

FNAME	M LNAME	SSN	BDATE	ADDRESS		S	SALARY	SUPERSSN	DNO	
John	B Smith	123456789	09-JAN-65	731 Fondren, Hous	ton, TX	м	30000	333445555	5	C^{1}
Franklin	T Wong	333445555	08-DEC-55	638 Voss, Houstor	, ТХ	М	40000	888665555	5	Site 3
Alicia	J Zelaya	999887777	19-JUL-68	3321 Castle, Spri	ing,TX	F	25000	987654321	4	Donort
Jennifer	S Wallace	987654321	20-JUN-41	291 Berry, Bellai	re, TX	F	43000	888665555	4	Depart
Ramesh	K Narayan	666884444	15-SEP-62	975 Fire Oak, Hum	ble, TX	M	38000	333445555	5	
Joyce	A English	453453453	31-JUL-72	5631 Rice, Housto	on, TX	F	25000	333445555	5	
Ahmad	V Jabbar	987987987	29-MAR-69	980 Dallas, Houst	on, TX	M	25000	987654321	4	
James	E Borg	888665555	10-NOV-37	450 Stone, Housto	on, TX	м	55000		1	
~					DNAME		DNUMB	ER MGRSSN	MGRSTARTD	TOTAL_S
Query at si	te 1:									
т Т					Research			5 333445555	22-MAY-88	
I _{Fname Lnan}	ne Dname (Emple	oyee $\bowtie_{Dno=}$	-Dnumher	Department)	Administrat	ion		4 987654321	01-JAN-95	
1 manie, Litan					Headquarter	s		1 123456789	19-JUN-81	
					Dummies			0 666884444	31-DEC-04	

 $F = \Pi_{Dnumber}(Employee)$ DNO $Q = \Pi_{Dname,Dnumber}(F \bowtie_{Dno=Dnumber})$ 5 4 DNUM DNAME Administration Research

COSC430 Lecture 5, 2020



 $\Pi_{Fname,Lname,Dname}(Q \bowtie_{Dno=Dnumber} Employee)$

	FNAME	LNAME	DNAME
er Department)	Franklin John	Wong Smith	Research Research
BER	Jennifer Alicia	Wallace Zelaya	Administratio Administratio
	Joyce	English	Research
4	Ramesh	Narayan	Research
5	Ahmad	Jabbar	Administratio



n n

n

19

Query or Update Decomposition

 $\Pi_{Fname,Lname,hours}\sigma_{dno=5}(Employee \bowtie_{Dno=Dnumber} Project \bowtie_{Pno=Pnumber} Works_on)$

 Suppose a site stores all information about projects working on these projects (Works_on5):

$$T_1 \leftarrow \Pi_{Essn,pno,hours} \sigma_{dno=5}(P$$

 $T_2 \leftarrow \Pi_{Fname,Lname,hours}(T_1 \Join_{Ssn=Essn} Employee)$

COSC430 Lecture 5, 2020

Idea: decompose query or update into sequence of queries or updates executed at the individual sites

controlled by department 5 (Projs5) and employees

Projs5 $\bowtie_{Pno=Pnumber}$ Works_on5)



Transactions

- A transaction is an atomic sequence of actions (reads and writes)
- ACID properties covered earlier
 - Atomicity
 - Consistency
 - Isolation
 - Durability

COSC430 Lecture 5, 2020

Success!

Begin() action action action action Commit()

Failure!

Begin() action action action action Rollback()



Transaction management

- Transaction T may touch many sites
 - $T = \{T_1, T_2, ..., T_n\}$
 - T_k runs at site k
 - How can atomicity be guaranteed?
- Two-Phase Commit protocol
 - Global transaction coordinator
 - e.g., the site that initiated the transaction can be the coordinator



Two-Phase Commit (2PC) protocol

- Request phase
 - Coordinator sends prepare message to all sites
 - Sites reply with "Ready" or "Abort" in their vote
- Commit phase
 - Coordinate sends a commit message if all sites vote "Ready", otherwise sends an abort message
 - Each site either commits or aborts, and replies with an acknowledge message





Three-Phase Commit (3PC) protocol

- and a site during commit phase
 - 3PC solves by introducing the Prepared to commit state (phase 2)
 - A site receiving preCommit knows that all other sites said Yes to canCommit?

2PC can't recover from failure of both the coordinator Coordinator Site canCommit? Phase 2 Yes preCommit Phase 2 ACK doCommit Phase 3

Coordinator





Site

Concurrency control in distributed DBs

- More complicated than that in a centralised DB environment
 - Dealing with multiple copies of the data item
 - Failure of individual site
 - Failure of communication links
 - Distributed commit
 - Distributed deadlock



Distributed locking

- Centralised approach
 - One dedicated site manages all locks
 - Bottleneck, unsalable, single point failure
- Primary-copy approach
 - Each data item has a primary site.
- Full distributed (voting) approach
 - A lock request is sent to all sites with a copy of the data item
 - for it
 - Use timeout to resolve deadlock

COSC430 Lecture 5, 2020



Each transaction asks the primary site for lock on the data item.

Each copy maintains its own lock and can grant or deny the request



Summary

- performance of a distributed database
- Queries in distributed databases can (sometimes) be optimised through the use of semi-joins
- Transaction management and concurrency control in distributed databases are more complex than the equivalent techniques within centralised databases

 Data fragmentation, allocation, and replication are key considerations that can significantly affect the

