

**MANAGING REPLICATION AND TRANSACTIONS
USING NEIGHBOUR REPLICATION ON DATA GRID**

NORAZIAH BINTI AHMAD

**Thesis Submitted in Fulfillment of the Requirement for the
Doctor of Philosophy in the Faculty of Science and Technology
Universiti Malaysia Terengganu**

July 2007

Abstract of thesis presented to the Senate of Universiti Malaysia Terengganu (UMT)
in fulfilment of the requirement for the degree of Doctor of Philosophy

**MANAGING REPLICATION AND TRANSACTIONS
USING NEIGHBOUR REPLICATION ON DATA GRID**

NORAZIAH BINTI AHMAD

July 2007

Chairperson : Professor Md. Yazid Mohd. Saman, Ph.D

**Member : Associate Professor Muhammad Suzuri Hitam, Ph.D
Professor Mustafa Mat Deris, Ph.D**

Faculty : Science and Technology

Replication is a useful technique for distributed database systems. Through this technique, a data object will be accessed (i.e., read and written) from multiple locations. Thus, it increases the data availability and accessibility to users despite site and communication failures. The all-data-to-all sites replication schemes such as Read-One-Write-All (ROWA) and Tree Quorum (TQ) are the popular techniques being used for replication and management of data in this domain. However, these techniques have its weaknesses in terms of data storage capacity and also data access times due to some number of sites must agree in common to execute certain transactions. In this study, the all-data-to-some sites scheme called Neighbour Replication on Grid (NRG) technique is proposed by considering only neighbours that have the replicated data. It is based on the logical structure of sites/servers in order to form a read or a write quorum in distributed database systems. The proposed technique considers only neighbours obtain a data copy. For simplicity, the neighbours are assigned with vote one and zero otherwise. The assignment provides a minimum communication cost with high system availability, due to the minimum

number of quorum size required. In addition, it minimizes the storage capacity as well as data access time.

A series of experiment was carried out by using three servers. Neighbour Replication on Grid (NRG) daemon is developed under Linux platform in the local area network (LAN) environment. It was carried out in Shell and Perl programming integrated with File Transfer Protocol (FTP) for the communications agent. The experimental results showed that the proposed model work successfully in managing replication and transaction when no failures occurred. Besides, the reconciliation and resolving conflict during system recovery are also supported when primary and neighbour replicas have failure.

ACKNOWLEDGEMENTS

I would like to express my most sincere gratitude to the supervisory committee Prof. Dr. Mustafa Mat Deris for his continuing support, professional guidance and for giving me an opportunity to learn what research is all about. Special gratitude also to Prof. Dr. Md. Yazid Mohd. Saman and Assoc. Prof. Dr. Muhammad Suzuri Hitam for their contributions, guidance and time towards this research. Not forget to Ms. Noraida Haji Ali, for her ideas during an early of this study.

Sincerely thanks should be forwarded to Vice Chancellor Universiti Malaysia Pahang (UMP), YH Prof. Dato' Dr. Mohamed Said bin Mat Lela for the KUKTEM Scholarship and also for the sponsorships to the International Symposium on DCABES 2006 and IIGSS Workshop 2007. Not forget to Dean of Faculty Computer System and Software Engineering (FSKPP), Prof. Dr. Abdullah Embong; Assoc. Prof. Ruzaini Abdullah Arshah; staffs in HR Department and all friends in UMP.

Special gratitude also to my family, especially to my father, Ahmad Zakaria; in memory of my mother, Esah Ibrahim; my beloved husband, Zakaria Mamat; my sister, Mazni Ahmad; my brothers Mohd Rashidi and Abdul Aziz; my mother in-law, Wan Sanah Wan Ibrahim; and my father in-law, Mamat Endut; for their patience and morale support.

Finally, I thank to all my friends especially to Rabiei, Norhayati, Nathrah, Suryani, Shuhadah, Ahmed, Che Norhaslida and others who have contribute in this research.

I certify that an Examination Committee has met on 8th of July 2007 to conduct the final examination of Noraziah Ahmad on her Doctor of Philosophy thesis entitled “Managing Replication and Transactions Using Neighbour Replication On Data Grid” in accordance with the regulations approved by the Senate of Universiti Malaysia Terengganu. The Committee recommends that the candidate be awarded the relevant degree. Members of the Examination Committee are as follows:

Mohd. Pouzi Hamzah, Ph.D.
Faculty of Science and Technology
Universiti Malaysia Terengganu
(Chairperson)

Md. Yazid Mohd. Saman, Ph.D.
Professor
Faculty of Science and Technology
Universiti Malaysia Terengganu
(Member)

Muhammad Suzuri Hitam, Ph.D.
Associate Professor
Faculty of Science and Technology
Universiti Malaysia Terengganu
(Member)

Mustafa Mat Deris, Ph.D.
Professor
Faculty of Information Technology and Multimedia
Universiti Tun Hussien Onn Malaysia
(Member)

Hamidah Ibrahim, Ph.D.
Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Independent Examiner)

Mashkuri Yaacob, Ph.D.
Professor
Universiti Tenaga Nasional
(Independent Examiner)

MOHD AZMI AMBAK, Ph.D.
Professor / Dean of Graduate School
Universiti Malaysia Terengganu

Date:

This thesis submitted to the Senate of Universiti Malaysia Terengganu and has been accepted as fulfilment of the requirements for the degree of Doctor of Philosophy.

MOHD AZMI AMBAK, Ph.D.
Professor / Dean of Graduate School
Universiti Malaysia Terengganu

Date:

TABLE OF CONTENTS

	Page
DEDICATION	ii
ABSTRACT	iii
ABSTRAK	v
ACKNOWLEDGEMENTS	vii
APPROVAL	viii
DECLARATION	x
LIST OF TABLES	xiv
LIST OF FIGURES	xv
LIST OF ABBREVIATIONS	xxi
 CHAPTER	
1 INTRODUCTION	1
1.1 Data Replication	2
1.2 Problem Statements	4
1.3 Objectives and Scopes of Research	6
1.4 Organization of Thesis	7
1.5 Conclusion	8
2 LITERATURE REVIEW	9
2.1 Read-One-Write-All (ROWA)	9
2.2 Voting (VT)	11
2.3 Tree Quorum (TQ)	13
2.4 Grid Structure (GS)	16
2.5 Three Dimensional Grid Structure (TDGS)	18
2.6 Conclusion	20
3 FUNDAMENTAL CONCEPTS AND THEORY	22
3.1 Formalization of Transaction	22
3.2 Serializability for Replicated Data	24
3.3 Quorum Intersection Property	29
3.4 Communication Costs	30
3.5 Availability	30
3.6 Conclusion	31
4 PERFORMANCE ANALYSIS OF EXISTING REPLICATION TECHNIQUES	32
4.1 Communication Cost Analysis	32
4.1.1 Read-One-Write-All (ROWA)	32
4.1.2 Voting (VT)	33
4.1.3 Tree Quorum (TQ)	34
4.1.4 Grid Structure (GS)	35
4.2 Operation Availability Model	36
4.3 Availability Analysis	37

4.3.1	Read-One-Write-All (ROWA)	38
4.3.2	Voting (VT)	39
4.3.3	Tree Quorum (TQ)	40
4.3.4	Grid Structure (GS)	41
4.4	Conclusion	42
5	NEIGHBOUR REPLICATION ON GRID MODEL	43
5.1	NRG Model	43
5.2	NRG Technique	44
5.3	The Correctness of NRG	50
5.4	Performance Analysis	51
5.4.1	Communication Costs	51
5.4.2	Availability	52
5.4.3	Optimal NRG	53
5.5	Conclusion	55
6	PERFORMANCE COMPARISONS WITH OTHER TECHNIQUES	57
6.1	Comparison of Communication Costs	57
6.2	Comparison of Availabilities	59
6.3	Conclusion	69
7	MANAGING TRANSACTIONS FOR NRG	70
7.1	Transaction Management	71
7.1.1	NRG Transaction Model	71
7.1.2	Failures Semantic	86
7.1.3	NRG Transaction Processing	88
7.1.4	NRG Transaction Manager (NTM)	92
7.2	Replica Management	96
7.2.1	The Coordinating Algorithm for the Primary Replica	96
7.2.2	The Cooperative Algorithm for Neighbour Replica	102
7.3	Reconciliation and Conflict Resolution	104
7.4	Correctness	109
7.5	Illustration Examples	112
7.6	Conclusion	121
8	IMPLEMENTATION OF MANAGING REPLICATION AND TRANSACTIONS	122
8.1	Hardware and software components	122
8.2	NRG Environment	124
8.3	An Application Example	126
8.3.1	NRG Daemon Configuration	128
8.3.2	Experiments and results	137
8.4	Conclusion	179
9	CONCLUSIONS AND FUTURE WORKS	180
9.1	Conclusions	180
9.2	Future Works	183

REFERENCES	186
APPENDIX	192
BIODATA OF THE AUTHOR	194

LIST OF TABLES

Table		Page
3.1	RD history H is a partial order with following condition	27
6.1	Comparison of the communication costs for read operation between ROWA, VT, TQ, GS, NRG under different set n of sites.	58
6.2	Comparison of the communication costs for write operation between ROWA, VT, TQ, GS, NRG under different set n of sites.	59
6.3	The read availability when $n = 36$ and $0.1 \leq p \leq 0.9$	61
6.4	The write availability when $n = 36$ and $0.1 \leq p \leq 0.9$	61
6.5	The read availability when $n = 64$ and $0.1 \leq p \leq 0.9$	64
6.6	The write availability when $n = 64$ and $0.1 \leq p \leq 0.9$	64
6.7	The system availability when $n = 36, p = 0.7$ and $0.1 \leq f \leq 0.9$	67
6.8	The system availability when $n = 64, p = 0.7$ and $0.1 \leq f \leq 0.9$	67
7.1	Meaning of Target Set	75
7.2	Primary-Neighbours Grid Coordination	78
7.3	The PNGC for $S(B_x) = \{8,3,7,9,13\}$	113
7.4	An example of how NRG handle concurrent transactions	116
8.1	Server main components specification	122
8.2	System development tools specification	124
8.3	The local IP address for each cluster member	125
8.4	Neighbour binary voting assignment for an experimental environment	126
8.5	Command input for the run levels	134
8.6	Primary-Neighbours Grid Coordination for replica A, B and D	137
8.7	Experiment result of how NRG handle concurrent transactions	145

LIST OF FIGURES

Figure		Page
2.1	A tree organization of 13 copies of a data object	14
2.2	An example of a write quorum required in TQ technique	15
2.3	A grid organization with 25 copies of a data object	16
2.4	An example of copies required to execute the read operation	17
2.5	An example of copies required to execute the write operation	17
2.6	A TDGS organization with 24 copies of an object	19
5.1	Examples of data replication in NRG	46
5.2	An assignment B for data file x where $S(B_x) = \{8,3,7,9,13\}$	49
6.1	The read availability when $n = 36$ and $0.1 \leq p \leq 0.9$	62
6.2	The write availability when $n = 36$ and $0.1 \leq p \leq 0.9$	62
6.3	The read availability when $n = 64$ and $0.1 \leq p \leq 0.9$	65
6.4	The write availability when $n = 64$ and $0.1 \leq p \leq 0.9$	65
6.5	The system availability when $n = 36, p = 0.7$ and $0.1 \leq f \leq 0.9$	68
6.6	The system availability when $n = 64, p = 0.7$ and $0.1 \leq f \leq 0.9$	68
7.1	Different sets of transactions requesting data x at sites i and j respectively ($B(i) = B(j) = 1$)	77
7.2a	Framework of semantics of NRG Transaction	81
7.2b	Initiate lock	82
7.2c	Propagate lock and obtain quorum	83
7.2d	Release lock, update and commit	84
7.2e	Failure (Unknown status)	85
7.2	Semantic of NRG Transaction	85

7.3	NRG Transaction Processing (for 4 clients)	91
7.4	An example of NRG transaction processing without system failure	113
7.5a	Primary replica 8 has system failures	117
7.5b	Neighbour replicas 8 and 13 have system failures	118
7.5	Different types of system failures while the transactions have performed	118
7.6	Replica priority check for neighbour replica 3, 7, 9 and 13	119
7.7	NRG_NTD entry for the first sequential T'_{γ_x, q_1}	119
7.8	NRG_NTD entry for the second sequential T'_{γ_x, q_1}	120
8.1	Data allocation on a cluster with 3 replication servers	125
8.2	NRG daemon is started while the system booting	127
8.3	NRG daemon is stopped when system halt	128
8.4	The “root” access is required to modify an NRG script	129
8.5	NRG script in an /etc/rc.d/init.d/ directory	130
8.6	The nrg_daemon script in the /usr/sbin/ directory	131
8.7a	The run level 0, 1 and 6 have links with an /etc/rc.d/init.d/NRG script	135
8.7b	The run level 2, 3 and 5 have links with an /etc/rc.d/init.d/NRG script	136
8.7	The particular run levels have links with an /etc/rc.d/init.d/NRG script.	136
8.8	NRG daemon has been successfully configured	136
8.9a	User “azie” requests to update the data file <i>dds</i> from replica A	138
8.9b	User “rosmawati” requests to update the data file <i>dds</i> from replica A	138
8.9c	User “noraziah” requests to update the data file <i>dds</i> from	138

	replica B	
8.9d	User “suryani” requests to update the data file <i>dds</i> from replica A	138
8.9	Users concurrently request to update the data file <i>dds</i>	138
8.10a	User’s information at the primary replica A	139
8.10b	User’s information at the primary replica B	139
8.10	User’s information at primary replica A and B	139
8.11a	$T_{\alpha a, q_1}$ gets and propagates the lock to its neighbours	140
8.11b	$T_{\alpha a, q_1}$ keeps propagating its lock in order to get a quorum	140
8.11	$T_{\alpha a, q_1}$ performs during an initialization and propagation lock phases	140
8.12	Kills <i>pid</i> of $T_{\alpha a, q_2}$ and broadcasts messages to user	141
8.13	$T_{\beta a, q_1}$ performs until user finishes update data.	141
8.14a	Primary replica B obtains the lock from neighbour replica D	143
8.14b	Primary replica B obtains the lock neighbour replica A	143
8.14	Primary replica B obtains the lock from neighbour replica D and A	143
8.15a	$T'_{\gamma a, q_1}$ change is committed at primary replica B	143
8.15b	$T'_{\gamma a, q_1}$ change is committed at neighbour replica D	144
8.15c	$T'_{\gamma a, q_1}$ change is committed at neighbour replica A	144
8.15	$T'_{\gamma a, q_1}$ change is committed at all replica of $S(B_a)$	144
8.16	User “noraziah” read the data file <i>dds</i> at replica B	148
8.17a	Primary replica B propagates its lock to the neighbour replica D and A	148

8.17b	Primary replica B obtains a quorum	149
8.17	A transaction performs during an initialization and propagation lock phases	149
8.18	The network error is simulated for the primary replica B	150
8.19a	Neighbour replica D detects the primary failure	151
8.19b	Neighbour replica A detects the primary failure	151
8.19	Neighbour replica D and A detect the primary failure	151
8.20	The content of the data file <i>dds</i> before the changes is made	152
8.21	Primary replica D propagates its lock to the neighbour replica A and B	153
8.22	User “noraziah” updates the data file <i>dds</i>	153
8.23	Primary replica D fails while committing data at neighbour replica B	154
8.24a	Neighbour replica A detects the primary failure	155
8.24b	Neighbour replica B detects the primary failure	155
8.24	Neighbour replica A and B detect the primary failure	155
8.25	The contents of the data file <i>dds</i> in Experiment 3	157
8.26	Primary replica A propagates the lock to its neighbours in Experiment 3	158
8.27a	Transaction change is committed at primary replica A	159
8.27b	Transaction change is committed at neighbour replica D	160
8.27	Transaction change is committed at all alive replicas of $S(B_{dds})$	160
8.28	The contents of the data file <i>dds</i> in Experiment 4	161
8.29	Primary replica A propagates the lock to its neighbour in Experiment 4	162
8.30	Transaction change is committed at primary replica A	162
8.31a	The contents of the data file <i>dds</i> for the first sequential T'_{γ_x, q_1}	164

8.31b	Primary replica B propagates its lock to neighbour replica D and A	165
8.31c	Primary replica B obtains a majority quorum	165
8.31d	Neighbour replica A checks its primary either is alive or not	165
8.31e	The network error is simulated for the primary failure	166
8.31f	Reconciliation for non-HPN replica	166
8.31g	Reconciliation for HPN replica	167
8.31	The first sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	167
8.32a	The contents of the data file <i>dds</i> for the second sequential $T'_{\gamma_{x,q_1}}$	167
8.32b	Primary replica A propagates the lock to its neighbours in Experiment 5	168
8.32c	User updates the data file <i>dds</i>	168
8.32d	Transaction change is committed at primary replica A	169
8.32e	Transaction change is committed at neighbour replica D	169
8.32	The second sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	169
8.33a	The contents of the data file <i>dds</i> for the third sequential $T'_{\gamma_{x,q_1}}$	170
8.33b	Neighbour replica B and D have failures.	170
8.33c	User updates the data file <i>dds</i>	170
8.33d	Commit and add new entries to NRG_NTD table	171
8.33	The third sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	171
8.34a	The contents of the data file <i>dds</i> for the fourth sequential $T'_{\gamma_{x,q_1}}$	172
8.34b	Neighbour replica D has been recovered from failure.	172
8.34c	Commit the fourth sequential $T'_{\gamma_{x,q_1}}$ and delete records in	172

NRG_NTD table for neighbour replica D

8.34d	Neighbour replica D catches the missed transaction and commits the fourth sequential $T'_{\gamma_{x,q_1}}$	173
8.34	The fourth sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	173
8.35a	The contents of the data file <i>dds</i> for the fifth sequential $T'_{\gamma_{x,q_1}}$	175
8.35b	Neighbour replica B and D have failures	175
8.35c	User updates the data file <i>dds</i>	175
8.35d	Commit and add new entries to NRG_NTD table	175
8.35	The fifth sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	175
8.36a	The contents of the data file <i>dds</i> for the sixth sequential $T'_{\gamma_{x,q_1}}$	175
8.36b	Neighbour replicas B and D have recovered from failure	176
8.36c	Primary replica A obtains all locks from neighbours	176
8.36d	Commit the sixth sequential $T'_{\gamma_{x,q_1}}$ and deletes all records of NRG_NTD table	177
8.36e	Catching all sequences of $T'_{\gamma_{x,q_1}}$ and commits the sixth sequential $T'_{\gamma_{x,q_1}}$	178
8.36f	Catching a sequence of $T'_{\gamma_{x,q_1}}$ and commits the sixth sequential $T'_{\gamma_{x,q_1}}$	179
8.36	The sixth sequential $T'_{\gamma_{x,q_1}}$ based on NRG Transaction Model	179

LIST OF ABBREVIATIONS

1SR	One-copy Serializable
CERN	European Organization for Nuclear Research
DDS	Distributed Database Systems
DDBMS	Distributed Database Management Systems
FTP	File Transfer Protocol
GS	Grid Structure
HEP	High Energy Physics
HPN	Highest Priority Neighbour
LAN	Local Area Network
LHC	Large Hadron Collider
NRG	Neighbour Replication on Grid
NRG_NTD	NRG <i>Need-To-Do</i> table
NTM	NRG Transaction Manager
P2P	Peer-to-peer
PNGC	Primary-Neighbours Grid Coordination
POSIX	Portable Operating Systems Interface
ROWA	Read-One-Write-All
SMTP	Simple Mail Transfer Protocol
TQ	Tree Quorum
VT	Voting
WAN	Wide Area Network

CHAPTER 1

INTRODUCTION

Several research articles have been published regarding distributed databases. Among them were those by Agrawal et al. [8, 9, 10, 22], Berstein et al. [46, 47], Chung [58], Garcia-Molina and Barbara [16], Maekawa [31], Mustafa et al. [33, 34, 35, 36, 66], Nicola [39], Stockinger [17] and Zhou et al. [68]. Those articles revealed that replicated data management is one of the current issues that still unsolved in distributed databases. Therefore, the study on this basis is initiated.

A database (DB) can be defined as a shared collection of logically related data that has been designed to meet the information needs of an organization and to be used by multiple users [15, 61]. The emergence of the network and additional communication facilities to a database system, can take it from centralize to a decentralize concept [11]. Distributed database system (DDS) is one of the major developments in the database area, where it moves from centralization that resulted in monolithic gigantic databases towards more decentralization [33, 35]. DDS is defined as a collection of multiple independent databases that operate on two or more computers that are connected and share data over the network [15, 61]. Meanwhile, a Distributed Database Management System (DDBMS) can be defined as the software that permits the management of the distributed database and make the data distribution transparent to users [61]. Nowadays, commercial database system such as *Oracle Database 10g* provides the required support for data distribution and inter-database communication

[53]. As a result of remarkable advance communication technologies, wireless and mobile computing concepts become reality. These concepts allow for even higher degrees of distributedness and flexibility in distributed databases [6, 35].

With the advances in distributed processing and distributed computing, the database research communities have carried out considerable works to address issues of data distribution, distributed design, distributed query processing, distributed transactions management and etc [61]. One of the major issues in data distribution is replicated data management. Typical replicated data management parameters are data availability and communication costs. These parameters share an inverse dynamic relationship: the higher the data availability with a low communication cost, the better the system is [31, 33].

1.1 Data Replication

Replication is an act of reproducing. It also addresses the management of the complete copying process [28]. Any type of data processing object can be implemented. These include the data files [13, 20, 23, 25, 42]; entire databases, specific tables, data within a specific tablespace [53]; service types such as Telnet, File Transfer Protocol (FTP), Simple Mail Transfer Protocol (SMTP) [38, 40] and etc. In a distributed database that relies on replication, the DDBMS may maintain a copy of fragment at several different sites [61]. Besides replicating, it also encompasses the administration to guarantee those data consistency across multiple sites.

Recently, scientific research and commercial application generate large amount of data that are required by users around the world. To illustrate this, the High Energy Physics (HEP) area deploys a new particle accelerator, the Large Hadron Collider (LHC) [17, 42]. It starts working at European Organization for Nuclear Research (CERN) in the year 2007. Several HEP experiments will produce Petabytes of data per year for decades. Those data will need to be managed and stored at more than hundreds of participating institutions. Thus, the data replication is a very useful technique to manage the large scale data across widely distributed communities.

Data replication plays an increasingly important role in this evolving world of distributed databases. Through this technique, an object will be accessed (i.e., read and written) from multiple locations such as a local area network or geographically distributed network world wide. For example, a student's results in a college, will be read and updated by lecturers of various departments. The financial instruments' prices will be read and updated from all over the world [36, 45]. Therefore, this technique provides high availability, fault tolerance and enhance the performance of the system [27, 33, 34, 36, 42].

Two approaches commonly used for replication, namely synchronous and asynchronous. Synchronous means move or operate together at the same time with each other while asynchronous is otherwise. Thus, synchronous replication provides what is called '*tight consistency*' between data stores. This means that the latency between data consistency is zero. If any copy is updated, the update immediately applied to all other copies within the same transaction. Data at all sites is always the same and exactly consistent, no matter from which replica the updated originated.

Conversely, asynchronous replication provides what is called '*loose consistency*' between data stores. This means that the latency between data consistency is always greater than zero. If one copy is updated, the change will be propagated and applied to the other copies within separate transactions. This copy changes may occur over seconds, minutes, hours, or even days later. Thus, some degree of lag always exists between the originating transaction that committed, and the effects of the transaction available at other replica(s) [28].

1.2 Problem Statements

Single centralized DB has low availability and reliability because if the DB site goes down the whole system fails. DDS has high availability and reliability. However, DDS introduces high redundancy as more than one site is used. This also creates low data consistency and data coherency as more than one replicated data need to be updated. Thus, some of the research problems that arise in DDS can be stated as follows:

- How does high availability of data achieved?
- How does redundancy reduced while there is an increase in the data storage capacity? How does replicated data minimized?
- How are synchronization mechanisms done in order to maintain the consistency of data when changes are made by transactions?

By storing multiple copies of data at several sites in the system, there is an increased data availability and accessibility to users despite site and communication failures. It is an important mechanism because it enables organizations to provide users with

access to current data where and when they need it. However, the storage capacity becomes an issue as multiple copies of data are replicated on different sites. Of course this way of data organization increases the data storage capacity. At the same time, expensive synchronization mechanisms are needed to maintain the consistency and integrity of data when changes are made by the transactions. This suggests that proper strategies are needed in managing replicated data in distributed database systems.

One of the simplest techniques is all-data-to-all sites replication scheme, namely Read-One-Write-All (ROWA) technique [55]. Read operations on a data object are allowed to read any copy, and write operations are required to write all copies of the data object. This technique has been proposed for managing data in mobile and peer-to-peer environment [6, 61]. It provides read operations with high degree of availability at low communication cost. However, it severely restricts the availability of write operations since they cannot be executed at the failure of any copy. This technique results in the imbalance of availability as well as the communication cost of read and write operations. The read operations have a high availability and low communication cost. Meanwhile, the write operations have a low availability with higher communication cost. Voting (VT) techniques [12] became popular because they are flexible and are easily implemented. This technique has been applied to the primary cluster for managing replicated data [20]. One weakness of these techniques is that writing an object is fairly expensive: a write quorum of copies must be larger than the majority of votes. To optimize the communication cost and the availability of both read and write operations, the quorum technique generalizes the ROWA technique. It imposes the intersection requirement between read and write operations [55]. Write operations can be made fault-tolerant since they do not need to access all

copies of data objects. Dynamic quorum techniques have also been proposed to further increase availability in replicated databases [21, 57]. However, these approaches do not address the issue of communication cost of read operations. Another technique, called Tree Quorum (TQ) technique [8, 9, 58] uses quorums that are obtained from a logical tree structure imposed on data copies. This technique has been proposed for persistent consistent distributed database commit in peer-to-peer network [3]. Nonetheless, this technique also has some drawbacks. If more than a majority of the copies in any level of the tree become unavailable, write operation cannot be executed. Several researchers have proposed logical structures on the set of copies in the database. To create intersecting quorums, the logical information has been deployed. The technique that uses a logical structure such as Grid Structure (GS) technique, executes operation with low communication costs while providing fault tolerance for both read and write operations [10]. However, this technique still requires that a bigger number of copies be made available to construct a quorum.

1.3 Objectives and Scopes of Research

The objectives of the research are as follows:

- i - To propose & to develop a new data replication model & to design its transaction management for the grid structure.
- ii – To analyze the model mathematically and to implement the model.
- iii – To test and to evaluate the performances of the proposed model.

This thesis only concentrates on the synchronous replication. Therefore, this research will focus on replication technique, in order to obtain high system availability with

low communication costs, in managing data replication by means of the synchronous replication.

In this research, the Neighbour Replication on Grid (NRG) technique is proposed based on the logical structure, in order to improve the read and write operations. Consequently, the implementation of this technique combines the replication and the transaction techniques. The motivations of this implementation are to show the clarity of the algorithms and also NRG technique can be used in the practical applications. A simplify approach to data distribution has been implemented by using the replication servers. A replication server handles the replication of data to remote sites. Through this mere simplification, it allows clearer presentation of an algorithm when less locking information of data is deployed.

1.4 Organization of Thesis

This thesis is organized as follows: Chapter 2 reviews the five major replica techniques, namely Read-One-Write-All, Lazy Replication, Voting, Tree Quorum and Grid Structure. In Chapter 3, the concepts of transaction, serializability, quorum intersection property, communication cost and availability are described. The performance analysis of the existing replication techniques are given in Chapter 4. Chapter 5 proposes the Neighbour Replication on Grid (NRG) technique. The performance comparisons with other replication techniques are then carried out in Chapter 6. Next, Chapter 7 presents the combination of replication and transaction techniques. In Chapter 8, the implementation of those techniques is described. Finally, the conclusion and the future works are presented in Chapter 9.

1.5 Conclusion

This chapter introduces the database systems and distributed database systems concepts. Data replication which is the main focus attention this research has been extensively elaborated. Together with advantages of data replication brings specific problems. It occurs since multiple copies of data are replicated at different sites. The storage capacity and expensive synchronization mechanisms to maintain the consistency and integrity of data become the issues. Thus suggests that proper strategies are required to solve the problems, which are the significance of this research.