

DISTRIBUTED DATABASES ON
LOCAL MULTIACCESS COMPUTER SYSTEMS

A Thesis

Submitted to the Faculty

of

Purdue University

by

Yao-Nan Lien

In Partial Fulfillment of the
Requirements for the Degree

of

Doctor of Philosophy

August 1985

ii

To my mother,
my brothers,
Shu-Man, and
Clinton

ACKNOWLEDGMENTS

I would like to express my sincere gratitude to my advisor, Professor Benjamin W.-S. Wah for his guidance, advice, and support during the research for this thesis. Without his inspiration, encouragement and friendship, this research could not be finished. I also wish to thank my committee members Professor Frederic J. Mowle, Professor Rangasami L. Kashyap and Professor Bharat Bhargava for their suggestions and comments. A special thanks goes to Professor Y. W. Ma of the Pennsylvania University for her precious comments. I also want to thank Mrs. Wah, the colleagues in our research group, and my friends who helped me to have a pleasant time during my study and research in Purdue University.

A special thanks goes to the Department of Computer and Information Science, The Ohio State University, for kindly providing me an excellent research environment since January 1986.

This thesis is fortunate to receive consistent support from the National Science Foundation under grants ECS80-16580 and ECS81-05981, and the GM-MAC, a research unit of Purdue University, sponsored by Purdue, Cincinnati Miliftron Corporation, Control Data Corporation, Cummins Engine Company, Ransburg Corporation, and TRW.

TABLE OF CONTENTS

	Page
LIST OF TABLES.....	viii
LIST OF FIGURES.....	ix
ABSTRACT.....	xiv
CHAPTER I - INTRODUCTION.....	1
1.1 Distributed Database vs. Centralized Database.....	2
1.2 Motivations of DDBMS's in Local Area Environment.....	4
1.3 Design Objectives and Design Issues of DDBMS's.....	6
1.3.1 Design Objectives.....	7
1.3.2 Design Issues.....	9
1.4 DDBMS's in Local Area Networks.....	15
1.4.1 Local Area Networks.....	16
1.4.2 Local Multiaress Bus Networks.....	21
1.4.3 End-to-end Transmission Speed of Multiaress Bus Networks.....	22
1.4.4 DDBMS's on Local Multiaress Bus Networks (DDBMN).....	23
1.5 Relational Data Model.....	24
1.6 Objectives of This Research.....	30
1.7 Overviews of Current Researches.....	31
1.8 Contributions of This Thesis.....	35
1.9 Significance of This Research.....	38
1.10 System Model.....	39
1.10.1 System Overview.....	39
1.10.2 Logical Database.....	39
1.10.3 Transactions.....	39
1.10.4 System Architecture.....	41
1.10.5 Architecture of Local Systems.....	42
1.10.6 Transaction Processing.....	45
1.10.6.1 Transactions and Processes.....	45

1.10.6.2 Transactions and Processes in DDBLMN	46
1.10.6.3 Transaction Processing Strategies in DDBLMN	48
1.11 Thesis Organization	51

CHAPTER III - FILE ALLOCATION

2.1 Introduction	52
2.1.1 Fragmentation	52
2.1.2 Fragment Allocation and File Allocation Problem	55
2.2 Simple File Allocation	58
2.2.1 Conventional SFAP	59
2.2.2 SFAP With Availability Constraints	60
2.2.3 SFAP With Average-delay Constraints	61
2.3 General FAP With Storage-Capacity Constraints	62
2.3.1 Optimal Solutions of FAP With Storage Constraints	64
2.3.2 Heuristic Solutions of FAP With Storage Constraints	67
2.4 Summary	67

CHAPTER III - DISTRIBUTED QUERY PROCESSING

3.1 Preprocessing Phase and Initial Broadcast	68
3.2 Introduction to Distributed Query Processing Problem	69
3.2.1 Redundant Materialization	70
3.2.2 Processing of Joins	73
3.2.2.1 Multi-site Joins and Semi-joins	73
3.2.2.2 Size Estimation and Selectivity	76
3.2.2.3 Static Semi-join Scheduling	77
3.2.2.4 Dynamic Semi-join Scheduling	78
3.3 Previous Work	79
3.4 Query Processing Strategy in DDBLMN	81
3.4.1 Five Phases in Query Processing	82
3.4.1.1 Concurrency-Control (CC) Phase	82
3.4.1.2 Local-Processing (LP) Phase	82
3.4.1.3 Global Semi-join (GSJ) Phase	83
3.4.1.4 Relations Transmission (RT) Phase	87
3.4.1.5 Post-Processing (PP) Phase	87
3.4.2 Improvements to The Five-Phase Query Processing	88
3.4.3 Protocol for Extremum Identification	90
3.4.4 Computing the Heuristic Value	94
3.5 Intelligent Query Scheduling	97
3.6 Intelligent Scheduling in DDBLMN	102
3.7 Summary	103

CHAPTER IV - CONCURRENCY CONTROL

4.1 Introduction	104
4.1.1 Consistency and Serializability	106
4.1.2 Throughput	106
4.2 Previous Work	107
4.2.1 Comparisons Among the Three Types Of Protocols	108
4.3 Concurrency Control in DDBLMN	110
4.3.1 General Description	110
4.3.2 The Locking/Unlocking Time	113
4.3.3 Minimum Locking Concurrency Control Protocol	115
4.4 Local Queries	122
4.5 Summary	122

CHAPTER V - GENERALIZED QUERY PROCESSING

5.1 DQP in Future Communication Networks	124
5.2 Introduction to GTSP	127
5.3 Solution Methods of GTSP	130
5.4 Graph Representation and Some Optimality Properties	132
5.5 TSP-Transformation Approach	137
5.5.1 G-tours on Graphs with Non-intersecting Clusters	137
5.5.2 G-tours on Graphs with Intersecting Clusters	144
5.6 Solutions For The Transformed Graph	146
5.7 Summary	167

CHAPTER VI - PERFORMANCE EVALUATION

6.1 Objectives of the Simulation	172
6.2 Queuing Network Model	174
6.3 Process Flow of DDBLMN	176
6.4 Assumptions in The Simulation	178
6.5 Simulator DDBLMN-SIM	180
6.6 Simulation Parameters	187
6.7 Simulation Results	191
6.7.1 Deadlock-free and Implementation Difficulty	194
6.7.2 Observations and Suggestions	197
6.8 Summary	210

CHAPTER VII - CONCLUSIONS.....	211
--------------------------------	-----

7.1 Summary and Conclusions.....	211
7.2 Suggestions For Future Research.....	214
7.2.1 High Bandwidth and Clustered Local Multiaccess Networks.....	214
7.2.2 Integrated Design of FAP and DQP.....	216
7.2.3 File Migration Problem.....	217
7.2.4 Intelligent Query Processing.....	217
7.2.5 Generalized Traveling Salesman Problem.....	218

LIST OF REFERENCES.....	220
-------------------------	-----

APPENDIX.....	232
---------------	-----

VITA.....	304
-----------	-----

LIST OF TABLES

Table	Page
1.1 Characteristic of local area networks.....	17
3.1 Relative performance of different heuristic schedules as compared to the optimal schedule.....	98
5.1 Performance of the TSP-transformation approach as solved by Karp's TSP algorithm.....	167
6.1 Parameters of DDBLMN-SIM system configuration.....	187
6.2 Configuration of relations in DDBLMN-SIM.....	187
6.3 Initial selectivities of all joining attributes in DDBLMN-SIM.....	188
6.4 Relation allocation in DDBLMN-SIM.....	188
6.5 Service time of servers in DDBLMN-SIM.....	188
6.6 Other timing parameters in DDBLMN-SIM.....	188
6.7 Data format in DDBLMN-SIM.....	189
6.8 Parameters used for query generation in DDBLMN-SIM.....	189
6.9 Tunable variables in DDBLMN-SIM.....	192
6.10 Simulation results of DDBLMN-SIM.....	198

LIST OF FIGURES

Figure	Page
1.1 A simple model of a DDBMS.....	10
1.2 Logical architecture of a distributed database system.....	11
1.3 The relationships of design issues in distributed databases.....	14
1.4 Topologies of local area networks (a) star (b) ring (c) bus (d) tree.....	19
1.5 An example of relational database (a) relations (b) relational schemas.....	26
1.6 Examples of operations in relational databases.....	28
1.7 An example of join operation.....	29
1.8 DDBMIN (a) From user's point of view (b) in the physical model.....	40
1.9 Logical architecture of local sites.....	43
1.10 Bernstein's transaction processing model.....	47
1.11 Example of two transactions being processed by DDBMIN.....	49
1.12 The processing sequence of a transaction with three queries Q_1 , Q_2 , Q_3 (a) the transaction; (b) the processing sequence.....	50
2.1 Example to illustrate the fragmentation problem.....	54
2.2 Example to illustrate the general file allocation problem.....	66
3.1 An example to illustrate the procedure of target data identification.....	71
3.2 An example to illustrate the procedure of join and semi-join (a) relation R and S; (b) sending relation S to join with R; (c) semi-join.....	75
3.3 The set of values of an attribute, its complement, and the initial complement.....	85
3.4 Distribution of relations and the result of the LP Phase.....	89
3.5 An example of the schedule generated by the basic query processing strategy.....	89
3.6 An example of the schedule generated by the improved query processing strategy.....	89
3.7 Protocol to identify the site with the maximum.....	91
3.8 The sequence of events in identifying the maximum.....	93
3.9 Intelligent semi-join scheduling model.....	100
4.1 Process for concurrency control with minimum locking time.....	116
4.2 Procedure for checking precedence relation between two transactions.....	117
4.3 Procedure for maintaining the precedence graph.....	117
4.4 An example to illustrate the concurrent execution of four transactions.....	120
5.1 An example of Algorithm DQPCTSP.....	126
5.2 The clustering presentation of Figure 5.1.....	128
5.3 Some g-tours on a graph.....	134

5.4	A g-tour T containing mutually exclusive nodes (v_x and v_y).....	136	5.21	Graph representation of Figure 5.2.....	165
5.5	A g-tour T' that dominates g-tour T in Fig. 5.4.....	136	5.22	I-N graph of the graph in Figure 5.21.....	166
5.6	Clustered graph G with non-intersecting clusters.....	139	5.23	G-S graph of the graph in Figure 5.22.....	168
5.7	Arcs added to G' in Steps (iv), (v) and (vi) of the construction process.....	140	5.24	ATSP tour found from the G-S graph.....	169
5.8	Proof of Theorem 2.....	142	5.25	The TSP tour on the I-N graph.....	170
5.9	An example of I-N transformation: (a) a part of graph G ; (b) the I-N graph G'' of G ; (c) a g-tour T'' on G'' ; (d) a better g-tour T''' on G''	146	5.26	Final G-tour on the original graph.....	170
5.10	A graph G and its I-N graph G''	150	6.1	Queuing model of simulator DDBLMN-SIM.....	175
5.11	A tour on G'' that visit v_x on G twice.....	150	6.2	Process-flow diagram of DDBLMN-SIM.....	177
5.12	A tour on G'' that does not visit all expanded nodes of v_x	151	6.3	Event flow of DDBLMN-SMPL.....	182
5.13	A tour on G'' that is a combination of Case (3) and (4).....	151	6.4	DDBLMN-SMPL implementation of DDBLMN-SIM.....	184
5.14	The new g-tour of T'' in Fig. 5.11 after canonical transformation.....	154	6.5	State diagram of DDBLMN-SIM.....	186
5.15	The new g-tour of T'' in Fig. 5.13 after canonical transformation.....	154	6.6	Testing results of DDBLMN-SIM (a) mean response time, (b) deviation of response time.....	191
5.16	Example graph 1: optimal and approximate GTSP solutions.....	160	6.7	Mean response time.....	200
5.17	Example graph 2: optimal and approximate GTSP solutions.....	161	6.8	Mean NETWORK queue length.....	201
5.18	Example graph 3: optimal and approximate GTSP solutions.....	162	6.9	Mean DISK queue length.....	202
5.19	Example graph 4: optimal and approximate GTSP solutions.....	163	6.10	Comparison of the utilization of DISK's and NETWORK.....	204
5.20	Example graph 5: optimal and approximate GTSP solutions.....	164	6.11	Comparison of the served time of CPU's, DISK's, and NETWORK.....	205
			6.12	The effect of packet size to the mean response time.....	207
			6.13	The effect of read_lock releasing time to the mean response time.....	208

6.14	The effect of materialization to the mean response time.....	209
------	--	-----

ABSTRACT

Lien, Yao-Nan, Ph.D., Purdue University, August 1986. Distributed Databases On Local Multiaccess Computer Systems. Major Professor: Benjamin W.-S. Wah.

Concurrency control, distribution design, and query processing are key issues in the design of distributed databases. In this research, we have studied these issues with respect to a relational database on a local computer system connected by a multiaccess broadcast bus. A broadcast bus not only allows information to be distributed efficiently but also provides a synchronization tool for the concurrency control. Consequently, it simplifies the distributed database design. Some NP-hard file placement problems are found to be solvable in polynomial time when updates are broadcast. A transaction model that integrates the concurrency control and query processing is proposed. In concurrency control, a protocol based on locking is used since much less communication overhead is needed for lock management. A dynamic strategy is used in query processing, as less data are transferred when compared to a static strategy. The status information needed in dynamic query processing can be conveniently obtained by broadcasting. Lastly, a new query processing strategy is proposed for the distributed database connected by the fast communication networks. The problem is formulated as a generalized traveling salesman problem, and is solved by various solution algorithms. The performance of the system is evaluated by simulations on a SUN workstation. The results shows the design is feasible and efficient.

CHAPTER I

INTRODUCTION

With digital computer technologies, Database Management Systems (DBMS's) provide better service for maintaining and accessing information over conventional information processing technologies. Expensive computer systems and the lack of computer networking technology in 1970's had limited the DBMS to be centralized in a computer complex, consisting of a few large computers, referred to as *centralized DBMS's* (CDBMS's). With the advance of computer networking technologies, a *Distributed Database Management System* (DDBMS) allows a single database to be efficiently and cost effectively spread over a computer network that consists of a number of geographically dispersed computers interconnected by a communication network. In the last decade, DDBMS's have been studied by a number of researchers [APE83] [BHA82a] [BHA83] [BER81a] [BER81b] [BER81c] [CER83] [CH184] [CHU82] [DOW82] [HEV79a] [PAP79] [RAM83] [SAC84] [WAH85] [WON77] [HEV85] [YU84a] [YU84b]. There are some pioneering systems developed: SDD-1 [ROT80], System R* [LIN84], Distributed INGRES [STO77], and LOCUS [POB83]. As the need of office automation and the capability of smaller computers grow rapidly, a DDBMS is now considered as a good candidate offering better information processing service for an environment in which users of a single database are spread within a limited distance area, referred to as *local area*. University environment is a typical example of such environments. DDBMS's in local area

environment are getting more and more attentions [CHA83] [HEV85] [KER82] [NGU81a] [SAC84] [WAH85]. In this research, we aim at studying the DDBMS's in local area environments and providing a feasible configuration accomplished with adequate control functions.

1.1 Distributed Database vs. Centralized Database

The major advantages of DDBMS's over CDBMS's are their superior achievements in reliability, data availability, throughput, and flexibility. Under some circumstances, a DDBMS may also provide faster response time and lower cost. On the other hand, the control functions in DDBMS's are much more complicated. There is a good discussion in Sykes's paper [SYK80].

- (1) **Reliability.** Reliability is the probability of a system keeping functioning in a distributed environment, failures on some components should not disable the entire system. Users' requests at the failed sites can be transferred to other sites; failures on a communication link can be handled by routing the traffic to alternate links. Although high reliability can also be achieved in centralized systems, the cost is much higher in general.
- (2) **Availability.** Availability of data is the probability of at least one copy of each piece of data being available in the system. In addition to the reliability, data availability is also a desirable feature of DBMS's. By replicating the data at different sites in DDBMS's, the chance that all access paths to a particular piece of data are broken is greatly reduced. The availabilities of other resources such as processing power may also be enhanced by the resources sharing capability of distributed systems.

(3) **Throughput.** Throughput is usually measured by the number of service requests (called *transactions*) that a system can process within a time unit. In a DDBMS, a transaction may be carried out by multiple processing units at different sites, or by the processing unit at any site with idle processing capacity. In other words, DDBMS's could have better load balancing, and hence larger throughput.

(4) **Flexibility/Extensibility.** The flexibility of upgrading system capacity is desirable whenever the future demand on the capacity is unpredictable at the design stage or is changing rapidly. Without such flexibility, a system may be over-specified at the design time to avoid the potential mismatch between the capacity and the future demand. Also, even if the future demand is predictable, the utilization of a system designed to match the future demand is very low before the demand matches the capacity. Those drawbacks can be avoided by allowing the system to be dynamically upgraded in response to changing demand. In a well-designed distributed system, adding new components or replacing old components have only little impact to the entire system. Consequently, high extensibility of DDBMS's is relatively easy to achieve, since the dependency among all components is usually much lower than that in CDBMS's.

(5) **Local Autonomy.** Each local site can have more design freedom to meet specific requirements of its local environment. Usually, the requirements of different user groups (eg., engineering, accounting, marketing, ... etc) in an organization may be very diverse and, even worse, may sometimes conflict with each other. In a DDBMS, each individual need can be easily satisfied

since each site can have its own autonomy. In contrast, it's very difficult for a CDBMS to satisfy all different needs. Even when tradeoffs among these needs are made, none of the users might be completely satisfied with the result.

(6) **Response time.** Most of the transactions in a well-designed DDBMS with good resource allocation strategies may have a great chance to be processed locally without extensive network activities. The average response time can be very short in this situation. However, certain types of transactions in DDBMS's may involve many complicated control functions making the response time relatively high comparing to CDBMS's.

(7) **Cost.** Accessing a centralized database from a remote site may involve expensive communications. A well-designed DDBMS may localize most of database accessing and hence, save expensive communication cost. However, under some circumstances, the overhead of the complicate control functions may be too high to be compensated by the benefit obtained from the localization of remote accesses.

1.2 Motivations of DDBMS's in Local Area Environment

In addition to the advantages described in Section 1.1, there are two driving forces that make the DDBMS a good candidate to provide good information processing service in a local area environment with physically dispersed users. The diameter of a local area environment is usually within a few kilometers. A hospital, a building with branch offices, or a manufacturing factory are some of the examples.

The first driving force comes from the need to improve office productivity due to the rapid growth of information volume and personnel cost. *Office Automation*, which tries to computerize most of the office functions, has been considered as the most promising approach to achieve this goal. *CAD/CAM/CAE*, as another example, have attracted a lot of attention for its potential to enhance engineering and manufacturing productivities. In these environments, information storage, maintenance, and retrieval are considered important services to support higher productivity.

The other driving force comes from the advance of technologies. Fast growing microelectronics and local area networking technologies in recent years have made local area networks, which connect smaller computers with reasonable processing power, inexpensive enough such that almost all companies can afford to install one.

The cost of fabricating a given logic element (processor, memory, peripheral controller) has been declining by a factor of 10 every three years over the last decade. Another way of looking at technological advances is to compare the size of memory on a single chip. Starting from 4 Kbits in 1973, it increased to 256 Kbits in 1979, and 1 Mbits DRAM today. The present density of logic on a chip is quite high such that any complex logic element can be built with a few chips at a cost lower than \$1000. A workstation with 1-Mbyte main memory, 1-MIPS (million instruction per second) processing power, and 1-Mbps networking capability is available in less than \$10,000. Current prices may fall in the future because of even larger production volumes and the rapidly developing trend to standardize chip interface logic. Further, smaller systems are much easier to

* CAD/CAM/CAE : computer aided design, computer aided manufacturing, and computer aided engineering.

operate and maintain than larger systems. This is even more important to small enterprises, since the shortage of skillful system engineers of large systems has made the central computing centers difficult to set up, even if they are affordable.

Recently, local area networks (LAN) have been extensively developed and used. These networks make use of such technologies as twisted pair, coaxial cable, radio transmission medium, as well as more sophisticated technologies like fiber optics. LAN can be installed in a building, on a campus, in a factory, etc. Many experimental and commercial LAN's have been in use for years, and both their feasibility and cost-effectiveness have been widely demonstrated.

Although CDBMS's may outplay DDBMS's on the efficiency in local area environments, a DDBMS is still considered as a better alternative for all its advantages other than efficiency and implementation complexity. In fact, the implementation complexity is not an important factor from the users' point of view as long as a good user interface is provided.

1.3 Design Objectives and Design Issues of DDBMS's

The objectives of DDBMS's are different from system to system. Most commonly addressed objectives are data integrity, cost effectiveness, good user interface, robustness, and simplicity. In fact, these objectives are generally applicable to most DBMS's.

1.3.1 Design Objectives

- (1) **Data integrity.** In a DBMS capable of processing more than one transaction at a time, transactions may interfere with each other such that either unexpected results are produced or the database itself is converted into incorrect states. Therefore, maintaining data integrity in processing transactions is one of the most important objectives of any DBMS.
- (2) **Cost effectiveness.** A DBMS should maximize its performance within a cost constraint, or in another way, meet the performance requirements with the least cost. The performance requirements are also different from system to system. These may range from minimizing response time or cost, to maximizing throughput or reliability, or a combination of more than one parameter.
- (3) **Good user interface.** Without a good user interface, a high performance system may be of little use or even useless to the users. In addition to good query languages as needed in CDBMS's, distribution transparency is also needed in DDBMS's [CER84]. Users do not have to know the location of any piece of data; the update to multiple copies of any piece of data should be automatically taken care of by the system itself. Further, a DDBMS may consist of some different local databases. It would be very difficult for users to access the database through different local systems a global data model is not available. Therefore, a global conceptual data model is desirable.
- (4) **Robustness.** This is the capability of a system to operate correctly in the presence of any failure or unexpected transactions. For example, communication-link failures may partition a system into isolated parts. A robust system should be able to continue the services with possibly degraded performance even if the hardware fails. Also, it should be able to

- recover the system back to a correct state, once the failed components are operational. It is generally true that the more complex a system is, the higher the probability that failures may occur. Therefore, the system designers should anticipate the failures and provide some ways to recover the system from failures automatically. In a worse case, a DBMS may operate with undetected errors if the failures do not stop the system. For example, a noise on the data path along which an update request is transmitted may change the database incorrectly if the system does not have error-detection and error-correction capabilities. As another example, two or more transactions may deadlock to each other. The system may eventually stop if it does not have any deadlock-prevention or deadlock-detection mechanisms. Robustness is particularly important in a DDBMS since its control functions are generally much more complicated than those in a CDBMS.
- (5) **Simplicity.** For a complicated system such as a DDBMS consisting of many components, it is difficult to prove its correctness, to make it robust, to document it, to maintain it, and to update it. Hence, simplicity is also an important objective in designing a DDBMS. However, very often, this can only be done at the cost of efficiency and performance. Compromises should be made in this case.

1.3.2 Design Issues

Although there is almost no difference between CDBMS's and DDBMS's on their objectives, there are some problems in DDBMS's that cannot be handled well by techniques developed for CDBMS due to the difference between their operating environments. The followings are some key design issues in DDBMS's.

- (1) **Architecture.** The design of a DDBMS is a complex and difficult task requiring careful consideration in all design aspects. A systematic architectural design can help designers solve the problem easily. In the simplest view, a DDBMS consists a set of computer systems interconnected by a communication network as shown in Figure 1.1. Each local system is equipped with a local database. According to the homogeneity of local systems, three types of integration can be used [PEE78]. In the first approach, each local DBMS 'knows about' the existence of others. All the components work together cooperatively like a single integrated DBMS system. This approach is possible only if the entire DDBMS is designed and built from the ground up. If this is not the case, existing DBMS's should be used. One of the two forms of 'federated' models should be used. Integrators and translators must be used to integrate different DBMS's into an integrated DDBMS system and to translate different database models into a unified conceptual model. The logical model of this approach is shown in Figure 1.2. More detailed discussions can be found in Peeble's paper [PEE78].

- (2) **Computer network.** This is one of the essential parts of the architecture design, providing the service for information exchange and resource

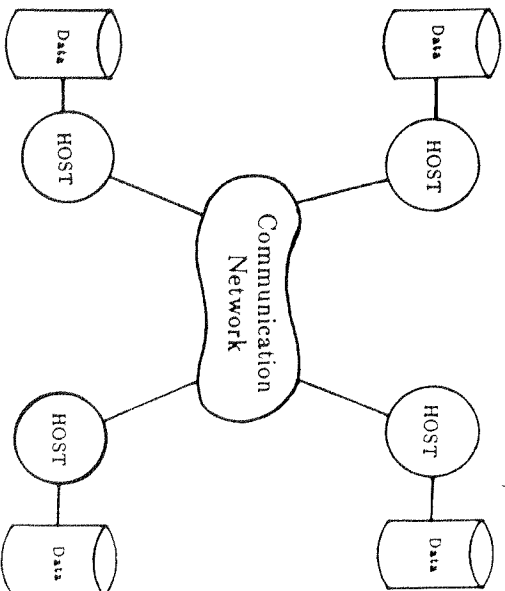


Figure 1.1

A simple model of a DDBMS.

ARCHITECTURE OF DDBMS

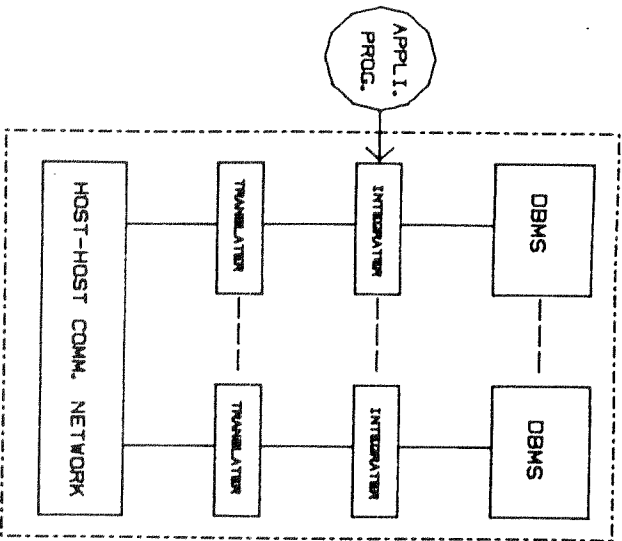


Figure 1.2

Logical Architecture of A Distributed Database System.

sharing. Basically, a network is a collection of two or more processing elements that offer various services and capabilities. A network consists of a certain type of transmission medium connecting all processing elements in a certain topology; the control mechanisms controlling the transmissions over the medium, some type of network interfaces, and a set of protocols controlling the transmission of information among the basic hardware elements. The combination of different choices on these components may result in a large number of networks differing in characteristics such as propagation delay time, communication capacity, addressing scheme, transmission cost, and effective network diameter. A good DDBMS needs a carefully designed computer network since most of the problems in DDBMS's cannot be solved well without considering the characteristics of the underlying networks.

(3) **Data Distribution.** The design of data distribution considers the allocation of data such that data can be accessed and maintained efficiently. In a DDBMS, replicated copies of a piece of data are allowed to be stored at different sites such that some expensive remote accesses can be localized. On the other hand, the overhead of maintaining the integrity of replicated data, the directory information, and the need for extra storage inhibits data from being over replicated. A careful allocation for the replicated data to balance the benefit and overhead is needed to maximize the cost effectiveness [CHU69] [DOW82].

(4) **Concurrency Control.** Concurrency Control maintains data integrity in a multiuser DBMS. It is complicated by some factors in DDBMS's as

follows. Multiple copies of a piece of data being stored, accessed, and updated at different sites at the same time; communication delays prohibit instantaneous distribution of status information such that it is difficult to maintain all sites in an identical state at all times; and the unexpected failures of some components may cause the database in different sites to be modified in different ways[BFR81b].

- (5) **Query Processing.** The operations performed on the data in a database is called a *query*. In a DDBMS, a query may originate from any site in the system. When processing a query, either the required files are assembled at a single site and the query is processed there, or the query and the intermediate results are sent and processed sequentially through the sites containing the files. A combination of the two strategies is also possible. The design of effective strategies of processing queries to maximize the system efficiency is known as the *DQP*, or *distributed query processing problem* [APE83] [YU84b].

Some other problems such as failure recovery, security/privacy, and logic database design are also relatively important.

All these issues are interrelated together as shown in Figure 1.3. For example, the files are partitioned and placed according to the characteristics of the network and the query; the design of the query processing strategy, the concurrency control algorithm, and the communication network depend on the placements and partitioning of the files. The issues shown in Figure 1.3 are also interrelated to other important problems, such as failure recovery, logical database design, and directory management. It is very difficult to solve these

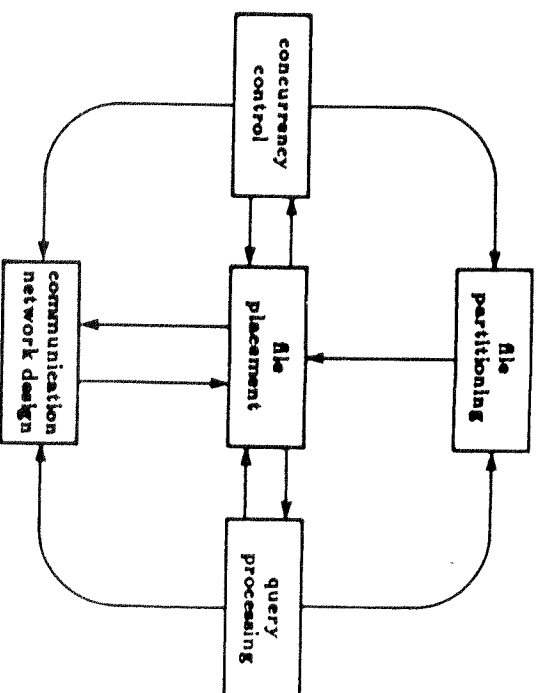


Figure 1.3 The relationships of design issues in distributed databases.

problems as a whole, so the designer usually decomposes them into independent problems and studies each based on simplified assumptions of others.

One potential problem of the above approach is that oversimplification may lead to undesirable results. An example is shown in the use of a network model that is independent of the hardware characteristic. A simple network model is often used because, besides simplifying the design of DDBMS strategies, it allows the software to be transportable, and permits many applications of different requirements to share the same network. A protocol hierarchy is used, which allows the DDBMS to interact with the lower levels through system calls. As a result, the status information needed in the control and the interchange of data among sites must be formulated into messages that are recognized by the lower levels of the hierarchy. This mismatch between the characteristic of the physical network and the requirements of the DDBMS results in inefficiency and increases the complexity of control strategies. To improve the performance, the capabilities of the network must be taken into account in the design. Based on this consideration, we have developed the objectives and the methodology of this research in Section 1.5.

1.4 DDBMS's in Local Area Networks

Since most of the problems of DDBMS's are so strongly dependent on the supporting network, it would be a good methodology to design the supporting network first and then solve the other problems based on the characteristics of the network. As a matter of fact, local-area-network (LAN) technology is getting matured and standardized * [IEEE83]. There is no need to design a special LAN * In 1980, the IEEE set up a Technical Committee, known as the 802 Committee, whose purpose is to establish standards for local area networks.

to support the DDBMS on local area environment if any existing LAN can support the DDBMS well. Therefore, the best methodology is to choose an existing LAN whose characteristic is close to what a DDBMS may need and to design the DDBMS based on the chosen LAN.

1.4.1 Local Area Networks

A *local network* is a communications network that provides for the interconnection of a variety of data-communication devices within a small area. A *local area network* (LAN), a special case of local networks, is a general purpose local network supporting minis, mainframes, terminals, and other peripherals [STA84]. The data rate can range from 0.1 to 20 Mbps, the distance from 0.1 to 50 kilometer, and error rate from 10^{-8} to 10^{-11} . The principal technological alternatives that determine the nature of a LAN are the *transmission medium* and the *topology*. There is a good survey in Stalling's paper [STA84].

- (1) **Transmission Medium.** Table 1.1 lists the types and the capacities of the transmission media that can be used in a LAN. The most commonly used media are twisted pair, coaxial cable, and optical fiber.
 - (a) **Twisted pair wiring.** Twisted pair wiring is the simplest one with the lowest capacity and cost. It is the most cost effective choice for single-building, low-traffic requirements. Only baseband transmission, which uses digital signaling, is suitable for twisted-pair wiring. One weakness of twisted-pair wiring is in its susceptibility to interference and noise.
 - (b) **Coaxial Cable.** High-performance requirements can be met by coaxial cables, which provide higher throughput, can support a large number of devices, and can span greater distances than twisted pair.

Table 1.1 Characteristic of local area networks.

Transmission medium	Twisted pair, coaxial cable, optical fiber
Topology	Bus, tree, ring, star
Transmission speed	0.1-20 Mbps
Maximum distance	~ 25 kilometers
Number of devices supported	10's-1000's
Attachment cost	500-5000

Both baseband and broadband transmission methods can be employed on a coaxial cable. In general, baseband transmission is simpler but less capable than broadband transmission. Baseband systems are typically from 1 to 10 Mbps and are generally limited to a single building. Broadband systems can support multiple data paths, each with 1 to 10 Mbps capacity. Existing CATV cabling technologies can be used in broadband LAN.

(c) **Optical Fiber.** Optical fiber has greatest capacity and is a promising candidate for future LAN's. It also has the following advantages over both twisted-pair wiring and coaxial cables: light weight, smaller diameter, low noise susceptibility, and no emissions. However, it has been little used so far due to its high cost and technical limitations. It is only suitable in to point-to-point configurations. The problem of insertion-tap loss must be overcome to use it in multipoint configurations.

(2) **Topology.** Characteristics of LAN's are also highly dependent on the topology used. Star, ring, bus, and tree are the most commonly used topologies. They are shown in Figure 1.4.

(a) **Star.** In star topology, a central switching elements is used to connect all the nodes in the network. Only one node can use the network to transmit or broadcast at any time. Therefore, a connection request to the central switching node is needed to establish a connection from one station to other stations.

(b) **Ring.** The ring topology consists of a closed loop, with each node attached to a repeating element, called *repeater*. Data circulate around the ring in a series of point-to-point data links between repeaters.

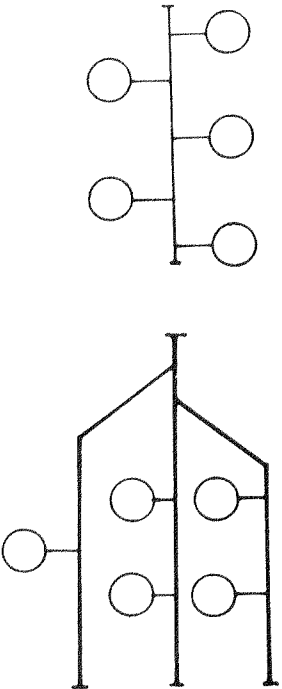
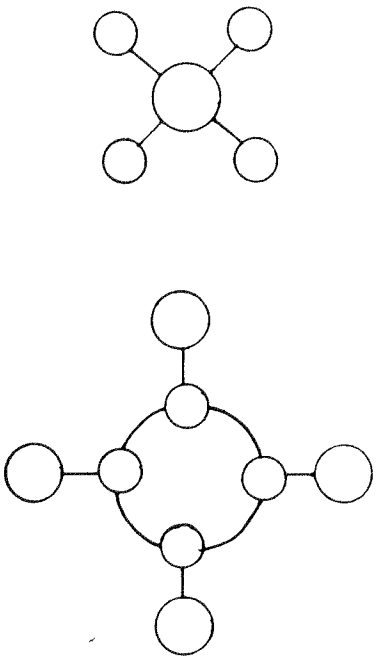


Figure 1.4 Topologies of local area networks (a) star, (b) ring (c) bus (d) tree.

Optical fiber is well suited in ring topologies such that its higher capacity can be utilized. Token rings, slotted rings, and insertion rings are the most commonly used access control mechanisms [STAK5]. A station wishing to transmit data should get the access right first and then transmit the data through the ring in packets with source and destination addresses. As the packet reaches the destination, the destination node copies the data. There are some problems associated with a ring topology. The reliability of the ring is jeopardized by a single failure on cables or repeaters. These failures are difficult to isolate. Further, adding new stations into the ring will degrade the performance.

- (c) **Bus/Tree.** Multipoint medium is used in bus or tree topology. The bus is a special case of the tree, in which there is only one trunk with no branches. Each station uses a passive interface to tap into the single bus. Since there is no active device on the bus, the reliability of bus LAN's is much higher than the ring topology LAN's. All stations share a common transmission medium, while only one station can use the medium to transmit at any time. A distributed-medium access protocol is needed to determine the station to transmit through the medium. Broadcast capability can be easily achieved since all stations 'listen' to the medium. The broadcast capability is very useful for DDBMS's. This will be discussed in the next section.

1.4.2 Local Multiaccess Bus Networks

Currently, the most popular bus LAN is the Ethernet [MET76] [SHO82]. Most of the low-cost, twisted-pair LAN's for microcomputers use a bus topology.

The following are the key characteristics of local multiaccess bus LAN's.

- (1) **Multipoint configuration.** More than one station are connected to the bus. Each transmitted message has the addresses of the source and destination stations. Bus access control protocols are needed.
- (2) **Serial transmission.** Since all stations share a single data bus, only one station can transmit at any time. This could become the bottleneck of the system. The utilization of the bus is an important issue in system design.
- (3) **Broadcast Capability.** All stations can listen to the bus simultaneously. The broadcast capability can be easily implemented. Many Ethernets support the broadcast capability [MET76] [SHO82]. Another useful potential capability is *multicasting*, in which messages may be targeted to more than one destination. This is also included in the Ethernet specification.
- (4) **Baseband/Broadband.** Both baseband and broadband transmissions can be employed except in the twisted-pair wiring bus, in which only baseband transmissions are allowed. If broadband transmission is used, more than one channel can be available. Virtually, it can be modeled as a multi-bus LAN in some cases.
As other multipoint configurations, a multiaccess bus needs bus access control. The access control is the protocol to determine which station has the right to transmit data when more than one station wants to transmit. Either centralized or distributed control can be used. A centralized control may have

the reliability problem since the entire network will fail if the control station fails. Furthermore, a heavy traffic to the control station may be induced by the access requests generated by the slave stations. The most commonly used distributed protocol is the *carrier sense multiple access with collision detection (CSMA/CD)* protocol. A station wishing to transmit should first sense the existence of the carrier signal on the bus to avoid destroying the current transmission. It can only transmit messages when the bus is idle (when no carrier signal presents on the bus). Otherwise, it should wait and keep sensing the bus until it is idle. A *collision* happens when more than one station transmit messages at the same time. Once a collision is detected, the transmitting stations should stop the transmission immediately and restart the procedure after a random period of time. Usually, a 'binary exponential backoff' algorithm is applied to determine the range of time within which the next attempt to retransmit will be taken.

1.4.3 End-to-end Transmission Speed of Multiaccess Bus Networks

The most important characteristic of a LAN in designing a DDBMS is the end-to-end communication speed, which is defined as the amount of information that can be transferred from the memory of one station in the original form to the memory of the other station in a ready-to-use form in a time unit. It is also called the *effective data rate*. We use these two terms interchangeably in this dissertation. It is much slower than the raw data rate on these networks because the overhead on the communication protocols, such as error detecting/correcting, medium-access control, flags, addressing, buffering, routing, packetization, flow control, sequencing, and synchronization, may reduce the effective data rate dramatically. For example, the typical end-to-end

transmission speed of a 10 Mbps Ethernet is only 1 Mbps. The TCP/IP protocols on a VAX-11/780 computer running 4.2BSD UNIX with a Ethernet may have only 0.5-0.75 Mbps end-to-end transmission speed.

1.4.4 DDBMS's on Local Multiaccess Bus Networks (DDBMN)

The broadcasting capability and multi-point configuration are the most important reasons that make the local multiaccess network a good choice for supporting a DDBMS in local area environments.

- (a) The update of multiple copies of a single piece of data can be made by one broadcast. This not only saves operating cost but also simplify the design.
- (b) The costs of remote accesses and updates are site independent. This can also simplifies the design.
- (c) System status is almost completely available to all stations by monitoring the activities on the bus. The status information exchange is, therefore, extremely small.
- (d) Every station gets the information on the bus almost instantaneously, since it is a multi-point configuration. For example, the minimum contention-slot time, which is at least twice the time to propagate a signal from one end of the bus to the other end, in Ethernet specification is 512 *bit-times* [SHO82]. The maximum propagation delay on the Ethernet, hence, will not exceed 256 bit-times. Comparing this to the minimum packet size of 576 bits, the propagation delay will not cause any problem. If we assume that the database state cannot be changed within a packet time, it's reasonably

* A bit-time is defined as the time that the transmission system needs to transmit a single bit. For example, in a 1-Mbps system, 1 bit time is 1 μ second.

safe to assume that the status exchange is instantaneous.

- (e) The multiaccess bus is a synchronization tool. When messages are not lost, all messages arrive at each site in the same order as they are sent. For example, if broadcast message m_1 is received at one site before a message m_2 , m_1 is received at all sites before m_2 . Therefore, the communication bus provides a synchronized communication environment [BAN79].

1.5 Relational Data Model

The *data model* of a database provides a logical view of the database to the users that can help users to access and manipulate the database without the knowledge of complicated physical organization of the database [COB70]. The two basic elements of a data model are:

- (a) a mathematical notation for expressing data and relationships among data elements, and
- (b) operations on the data that serve to express queries and other manipulations of the data elements[ULL82].

Among three different models (relational, hierarchical, and network) that have been extensively studied, the relational model is the most popular one due to its simplicity, symmetry, easy to use, and its strong theoretical foundation. For the same reason, this research is based on the relational data model.

The mathematical concept underlying the relational model is the set-theoretic relation, which is a subset of the Cartesian product of a list of domains. A *domain* is a set of values. The *Cartesian product* of domains D_1, D_2, \dots, D_k , written as $D_1 \times D_2 \times \dots \times D_k$, is the set of all k -tuples (v_1, v_2, \dots, v_k) such that v_1 is in D_1 , v_2 is in D_2 , and so on.

A *relation* is any subset of the Cartesian product of one or more domains. It is always assumed finite. The members of a relation are called *tuples*. Each relation that is a subset of $D_1 \times D_2 \times \dots \times D_k$ is said to have *arity* or *degree* k . A relation can be viewed as a table, in which each row is a tuple. Columns are called *attributes*. The set of attribute names for a relation is called the *relation scheme*. If we name a relation R , and its relation scheme has attributes $a, b, c,$ and d , then the relation scheme is written as $R(a, b, c, d)$. The attribute a of relation R is denoted as $R.a$. The order of both the tuples and the attributes are unimportant. A couple of examples are shown in Figure 1.5.

An operation to a database can be decomposed into two parts, the target-data identification and the target-data retrieval or update. The target-data identification usually needs a sequence of basic operations to be applied on the database when users query the database using high level query languages. There is no unique way to process such queries. Also, the processing overhead for different processing sequences may be quite different. Hence, a major task of database design is to design a strategy that can produce efficient processing sequences for all possible queries.

There are five basic operations that can manipulate a relational database.

- (1) **Selection.** Selection selects those tuples in a relation that satisfy a given condition, and is denoted as $\sigma_F(R)$, where F is the condition and R is the relation to be operated.
- (2) **Projection.** Projection takes unique values in an attribute of a relation, and is denoted as $\pi_A(R)$, where A is the attribute to be projected and R is the relation to be operated.
- (3) **Join.** When two relations share a common attribute, they can be joined.

The result of a join operation on a common attribute of two relations is the

AJAX, INC.	NEW YORK
BROWN & SONS	LOS ANGELES
CARLTON ENTERPRISES	CHICAGO

(a)

COMPANY	
COMPANY NAME	CITY

(b)

AJAX, INC.	WIDGET	4.00
AJAX, INC.	GISMO	4.50
BROWN & SONS	DOJIGGER	5.00
(NULL)	WHATCHET	3.50

PRODUCT		
COMPANY NAME	PRODUCT NAME	COST

Figure 1.5 An example of relational database (a) relations (b) relational schemas.

Cartesian product of the two relations, with the tuples which share the same value in the common attributes. Only one common attribute is kept in the resulting joined relation. This type of join operations is also called equi-join. There are other types of join operations [UL1,82].

- (4) **Union.** The union of relations R and S is the set of tuples that are in R or S or both. It is denoted as $R \cup S$.
- (5) **Set difference.** The difference of relations R and S, denoted as $R - S$, is the set of tuples that are in R but not in S.

The examples of these operations are shown in Figure 1.6. These operations can be concatenated together to form more complicated queries. There is no unique way to form a complicated query to get the final result. Furthermore, different basic-operator sets may have equivalent manipulating power. That is, the result obtained by applying a sequence of operations belonging to an operator set can also be obtained by applying operations belonging to another operator set. For example, the join operation can be replaced by the Cartesian product operation, since a join operation can be replaced by a combination of Cartesian product, selection, and projection. This is shown in Figure 1.7. Notice

- (a) that the join operation is a binary operation with two relations as operands,
- (b) that the join operation is symmetric,
- (c) that the size of a joined relation would be much larger than the sizes of each individual relations, and

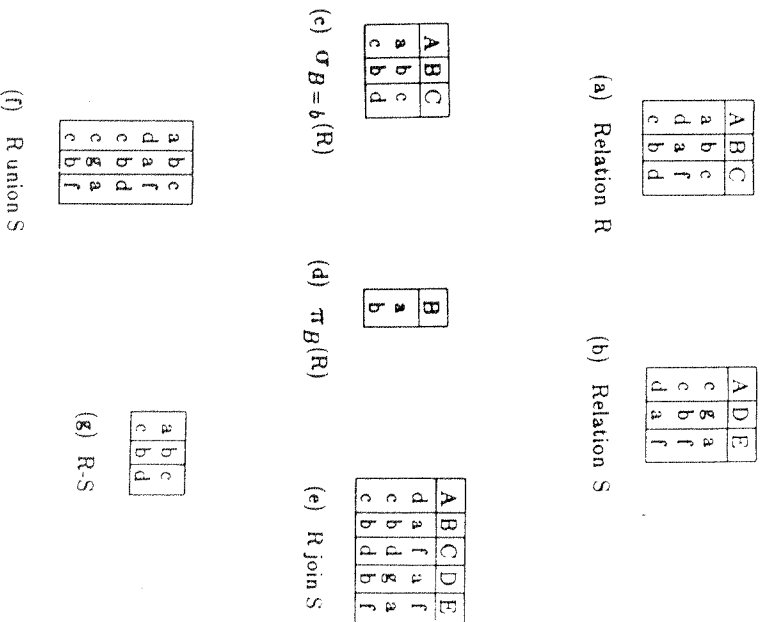


Figure 1.6 Examples of operations in relational databases.

* **Cartesian product.** Let R and S be relations of arity k_1 and k_2 , respectively. Then $R \times S$ is the set of $k_1 + k_2$ tuples whose first k_1 components form a tuple in R and whose last k_2 components form a tuple in S.

R		R Join S		S	
A	B	B	B	C	C
1	b			a	4
2	b			b	5
3	c			b	6

1. $R \times S$

A	B	B	C
1	b	a	4
1	b	b	5
1	b	b	6
2	b	a	4
2	b	b	5
2	b	b	6
3	c	a	4
3	c	b	5
3	c	b	6

2. Selection ($R.B=S.B$)

A	B	B	C
1	b	b	5
1	b	b	6
2	b	b	5
2	b	b	6

3. Projection A, B, C

A	B	C
1	b	5
1	b	6
2	b	5
2	b	6

Figure 1.7 An example of join operation.

(d) that there is a selection embedded in a join operation.

The most commonly used operations are selections, projections, and joins. Selection and projection are unary operations and generate relations with smaller size. A lot of overhead may be involved in a join operation since it usually expands the size of the joined relation. Further, in a DDBMS, to perform a join operation on two relations stored in two different locations may involve a lot of network activities. Therefore, to process join operations efficiently is the most critical task in the design of either a CDBMS or a DDBMS.

Most of the existing query languages are based on three primitive query languages: relational algebra, tuple relational calculus, and domain relational calculus [ULL82]. For simplicity, we consider the join queries only in this thesis, which are in the form of:

GET R.a WHERE (R.a = S.a) AND (R.b = T.b)

The result of the query consists attribute a only, which is projected from the result of joining relations R, S, and T together.

1.6 Objectives of This Research

This dissertation intends

- (1) to study the key problems of distributed databases in local multiaccess networks (DDBLMN) including distribution design, query processing, and concurrency control, as well as
- (2) to propose strategies for solving these problems in an integrated fashion, such that a DDBLMN supporting efficient execution of user transactions can be easily implemented.

We have mentioned that solving all the DDBMS problems together in an integrated fashion is very difficult. Further, there is no need to design the communication system with the DDBMS together. Therefore the local multiaccess network, which has nice characteristics that can simplify the design, is chosen as the supporting network. Then, the DDBMS is designed based on these characteristics.

To further simplify the problem, the distribution design problem is solved separately although it is related to the other two problems. Distribution design, is done before a DDBMS is operational while query processing and concurrency control problems is to be solved dynamically. As long as the system characteristics do not change rapidly, this approach can lead to stable and efficient design.

The relationship between query processing and concurrency control problems is so strong that it is difficult to combine the solutions of solving each problem individually.

1.7 Overviews of Current Researches

Most current researches of DDBMS's in local area environments have attempted to solve one or more individual problems. Actually, most of them consider the query processing problem without considering the concurrency control problem.

Hewner, Wu, and Yao [HEV85]

This research investigates distributed query processing on *address ring* and *broadcast networks*. Only communications cost is concerned.

An address ring is really a token ring [STA84], which uses a token to control the ring-access right. A station wishing to transmit a message to one or more other stations should get the token first. The addressees copy the messages down to the stations by matching their addresses to the addresses associated with the messages. The average transmission cost for a message is taken as the average time for the sending site to gain control of the token ring and a term proportional to the data volume and the distance to the addressees.

A broadcast network under consideration can have either a bus or ring topology with various protocols. The average communication cost of sending a message on the network is the average time for the sending site to gain control of the network and a term proportional to the data volume.

Static query optimization algorithms, in which the processing strategy is determined before a query is processed and cannot be changed during processing, are proposed. For some restricted cases, these strategies may be capable of getting optimum solutions. These strategies can be extended into heuristic algorithms to find strategies for general queries.

Sacco [SAC84]

Sacco investigates the query processing problem on broadcast networks. Three heuristic algorithms are proposed. These algorithms can be implemented either statically or dynamically. However, no adequate approach for implementing dynamic query processing is shown.

LAMBDA Systems [CHA83] [CHA84a] [CHA85]

LAMBDA, a database system on a local multiaccess/broadcast network, is an experimental system designed at the AT&T Bell Laboratories, Murray Hill, for studying DBMS issues in a local area network. It provides (1) high data availability through replicated data, (2) automatic recovery from site failures, and (3) uninterrupted transaction processing during site recovery. The design is simplified by assuming a reliable broadcast network. Algorithms for transaction commitment, concurrency control, and crash recovery are designed in LAMBDA.

Gouda and Dayal [GOU81]

Gouda and Dayal investigate the query processing problem in a local area network with one or more parallel broadcast buses. The database is assumed without data redundancy. This research pointed out that the query processing problem in this environment is *NP-hard* and heuristic algorithms are needed to find good solutions to the query processing problem.

MICROBE Systems [NGU81a] [NGU81b]

MICROBE is a DBMS on a network of LSI-II micro-processors connected by a local broadcast network implemented at the Laboratoire IMAG in the University of Grenoble. A mix of static and dynamic query processing strategies is used in MICROBE. The primary optimization objective is to minimize local processing and communication costs.

* A problem is called NP-hard if there is no known algorithm that can solve the problem in a time polynomially proportional to the problem size.

LU [LU85]

Lu proposes a *load-balanced query processing* (LBQP) approach to solve the query processing problem using load balancing in a local DBMS. The database is assumed to be fully replicated. A dynamic query processing strategy, which is capable of migrating a process to other sites during the query processing, is used in this research to enhance system load balancing and, consequently, reduces the response time.

UNITY Systems [KER82]

The UNITY system, developed at the AT&T Bell Laboratories, consists of a minicomputer and a set of microcomputers connected by a star network. Kerschberg, Ting, and Yao presented several query optimization techniques that are specially designed for star networks.

Ceri, Paolini, Pelagatti and Schreiber [CER82]

Ceri et al. were implementing a relational DBMS on a star local area network in Italy. The major concern of this system is on the file partitioning problem (see Section 2.1.1). Conventional majority two-phase locking concurrency control protocol (see Chapter IV) is used for concurrency control.

DDBLMN

Similar to LAMBDA system, DDBLMN in this research is assumed supported by a reliable local multiaccess/broadcast network. However, file allocation problem, distributed query processing problem, and concurrency

control problem are solved together in DDBLMN. While all existing systems only focus on one particular problem. The performance of such systems may be quite different from the expectation when they are really implemented in a working environment. Further, the design of DDBLMN takes the advantages of supporting network such that simple but efficient control strategies can be developed to enhance the system performance. The overview of working environment of DDBLMN is shown in the Section 1.10.

1.8 Contributions of This Thesis

The contributions of this thesis are as follows:

1. *A new transaction processing model that is particularly suitable for DDBMS on local multiaccess networks.*
The information exchange in general networks is inefficient such that a transaction is generally processed at the site it was entered. Either the needed information is gathered from other sites, or the entire transaction is migrated to related sites to be processed. These processing models do not take advantage of the distributed processing system, i.e., the parallelism is low. The proposed model allows all related sites to work together to carry out a transaction and to maintain a high degree of parallelism.
2. *Solution of the simple file allocation problem in local multiaccess networks in polynomial time.*
In general networks, SFAP is an NP-hard problem. In local multiaccess networks, we have solved in polynomial time the standard SFAP and the SFAP with availability constraints. Although the problem is NP-hard when there is an average-delay constraint, it is isomorphic to the knapsack problem, which has good approximate solutions.

3. *Efficient solution of the general file allocation problem in local multiaccess networks.*
The FAP with storage constraints in local multiaccess networks is NP-hard. An efficient algorithm combining a branch-and-bound search with the knapsack problem has been developed.
4. *Integration of distributed query processing and concurrency control.*
Currently, distributed query processing and concurrence control are studied independently in general networks. The combined results may not be good for an integrated system. An integrated model that combines query processing optimization with concurrency control has been developed in this thesis.
5. *Generalization of the distributed extremum identification algorithm on local multiaccess networks to identify the sub-query to be processed.*
Based on Wah and Juang's distributed extremum identification algorithm [WAH85b], which can identify the site with the extremum value efficiently, a general algorithm has been developed to synchronize transaction processings.
6. *Realization of dynamic query processing on local multiaccess networks.*
The major obstacle to realize dynamic query processing, which performs better than static query processing, lies in the difficulty to collect global system status in real time. We have proposed an efficient dynamic query processing algorithm based on the extremum identification algorithm.
7. *Performance improvement by redundant materialization and non-profitable data identification.*

Redundant data are used in query processing to help reduce the size of intermediate results. The non-profitable data identification is included in the query processing algorithm to eliminate useless transmissions. These two methods work together to improve the performance of query processing.

8. *Transformation of the query processing problem in future computer networks into the generalized traveling salesman problem.*

Current solutions for the query processing problem are not adequate for DDBMS in future computer networks, which are characterized by high-speed data links. We have modeled the problem as a generalized traveling salesman problem. An efficient algorithm that utilizes the existing solutions of the standard traveling salesman problem has been developed and analyzed.

9. *Development of a framework of intelligent query processing system.*

Knowledge based on the history of the database operation can be used to generate better query processing strategies. This provides a good direction to improve the query processing strategies.

10. *Development of an efficient concurrency control protocol in local multiaccess networks.*

Current concurrency control protocols are designed for general networks, which do not support efficient information exchange among local processing sites. Information exchanges in the local multiaccess networks are extremely efficient, and an efficient and highly concurrent protocol for DDBLMN has been developed.

1.9 Significance of This Research

The current trend in multi-user local-computer-system development tends to be the network of small computers. A number of interesting projects are in progress today on designing such systems. Among all applications in a local computer system, database management is definitely very important. As a consequence, the design of a robust, easy-to-use and cost-effective DDBMS in a local area environment is becoming an important issue. Although DDBMS on *wide area networks* have been studied extensively, their results cannot be applied directly for the following reasons. First, the characteristics of local area networks are quite different from global networks. Second, the users' behavior and users' requirements are different, e.g. the mobility of users in such an environment may be much higher. Lastly, the combination of solutions derived from individual problems in global networks may not perform well in an integrated local system. In short, a DDBMS study in an integrated fashion is needed for the local area environment.

Although the research in this thesis does not solve all the problems of DDBMS, we have studied three of the most important problems, namely, file allocation problem, distributed query processing problem, and concurrency control problem, in an integrated fashion. The feasibility and the effectiveness of the design are demonstrated by simulation studies. Based on this research, a prototype can be implemented easily. This research also provides a framework for DDBMS design in a local environment.

1.10 System Model

The DDBLMN is designed is based on the environment described in this section.

1.10.1 System Overview

From the users' point of view, the system is nothing but an integrated relational database consisting of relations as shown in Figure 1.8(a). Physically, it's a set of local database systems interconnected by a multiaccess/broadcast communication bus as shown in Figure 1.8(b).

1.10.2 Logical Database

It is assumed that a very large relational database with a global conceptual schema is available to the users. Virtually, users can access the database from different locations through the same conceptual schema. Therefore, we assume that the probability of a user accessing the database from any location is non-zero and is not uniformly distributed. There are some locations that a user may access more frequently than other sites. These two assumptions are essential to the distributed database. With a strong locality of access, there is no need to design an integrated database. On the other hand, a random access pattern is also unrealistic. Each relation is stored as a file and is the unit to be locked, i.e. a relation is either available or not available to a user, without partial availability. The statistics about the frequencies of access and update is only available on a relation basis. The data directory is fully distributed over all sites. The directory is maintained in the same way as a relation. Using fully replicated directories simplifies the system design dramatically. The additional storage is not significant since the directory is much smaller than the database itself.

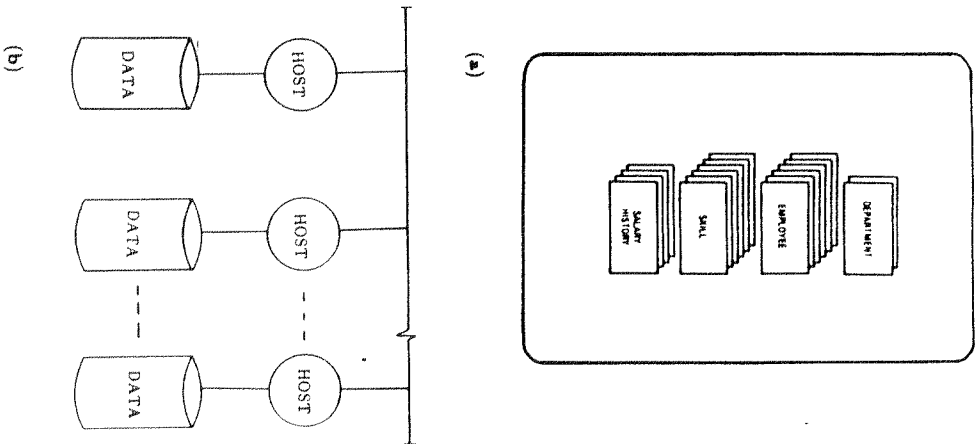


Figure 1.8 DDBLMN (a) from user's point of view (b) in the physical model.

Unless the database is reorganized very frequently, the update frequency to the directory is much less than the retrieval frequency. Moreover, the updating cost is independent of the number of copies in a network with broadcast capability. Therefore, the additional overhead of directory maintenance is also insignificant.

1.10.3 Transactions

The users activities to the system are called *transactions*. Each transaction is atomic, which should be either completely finished or completely aborted, leaving the system unchanged. All transactions are independent of each other. We further assume that each transaction can not be aborted by the user. Therefore, a transaction is always terminated successfully in a reliable system unless the concurrency control protocol aborts the transaction. A transaction may consist of a set of independent queries. Each query consists of two parts: the specification of the target-data and the operations on the target-data. The specification is in the form of consecutive equi-joins. An example is shown in the following.

```
GET      R.a  WHERE (R.a = S.a) AND (R.b = T.b)
UPDATE  S.a  WHERE (S.a = T.a) AND (R.b = S.b)
GET      T.a  WHERE (T.a = S.a)
```

Notice that the values to be updated are not given here since it is not important.

1.10.4 System Architecture

As shown in Figure 1.8(a), the DDBLMN consists of a set of local systems connected by a local multiaccess bus with the broadcast capability. A local database system is accomplished by each local system. The transmission system is assumed reliable, which simplifies the system design [CHA84a] [CHA84b]. The

transmission system provides guaranteed delivery of messages. A distributed multiaccess protocol, such as the CSMA/CD protocol, is assumed to control the bus accesses. In general, each station has equal priority to access the bus. However, the performance can be enhanced if priority messages are allowed. The end-to-end communication speed of the bus is about 1 Mbps (million bits per second). The number of local systems is around twenty to several hundred, each one with about 1 MIPS (million instructions per second) processing capability. The transmission cost, which can be either real cost or the elapsed time depending on the system objectives, is proportional to the volume of transmitted data. Under these assumptions, the processing capability per station is much larger than the communication bandwidth per station. Consequently, the local processing cost can be ignored.

1.10.5 Architecture of Local Systems

DDBLMN is not designed for any system with a specific local architecture. However, a local system architecture described here serves as an example. The logical structure of the local system at each site is shown in Figure 1.9. Each system consists of five different modules: Transaction Monitor, Transaction Handlers, File-Server subsystem, Concurrency-Control subsystem, and Network-Interface subsystem.

- (1) **Transaction Monitor (TM)**: This module coordinates all local and remote transactions by
 - (a) initiating/terminating Transaction Handlers that handle each individual transaction, and
 - (b) providing interface functions among all other modules.

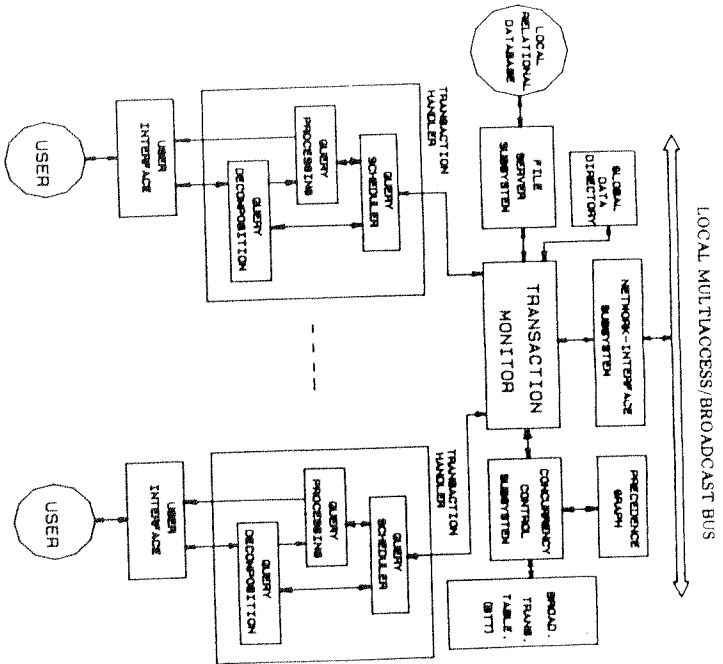


Figure 1.9 Logical architecture of local sites.

It also maintains a global directory and provides access transparency to users.

users.

- (2) **Transaction Handler (TH):** This module is created by the Transaction Monitor for handling each individual transaction. The TH can be initiated either locally or remotely. All the transaction processing functions, such as query decomposition, user interface, query processing, and query scheduling for a transaction, are performed in the TH designated for the transaction. It is eliminated when this transaction terminated.

- (3) **File-Server Subsystem:** This module, which serves as a local DBMS, is responsible for accessing and updating the local database. A TH makes relation access/update requests through the TM to the file-server subsystem.

- (4) **Concurrency-Control Subsystem:** This module monitors the actions of all local and remote transactions on the network through the Network-Interface Subsystem. All necessary status information is also maintained here. All the local data access/update requests should be checked by this module to maintain data integrity.

- (5) **Network-Interface Subsystem:** This module, which performs all the interface functions to the communication network, provides guaranteed end-to-end message delivery. The network interface functions include addressing, network access control, flow control, error control, message receiving, format conversion, protocol control, sequencing, synchronization, fragmentation, and reassembling.

The detailed relationship between TM and the TH's is discussed in next section.

1.10.6 Transaction Processing

Users of a DDBMS interact with the database by executing transactions. A transaction, which consists of a sequence of operations (queries) on one or more database objects, transforms a current consistent database state into a new consistent state. Query processing and concurrency control are two key tasks in transaction processing. In this section, we discuss how a transaction is processed.

1.10.6.1 Transactions and Processes

A process is a program initiated by the system for handling a particular task. A process is called *active* if it is currently executed by the system. In a concurrent environment, more than one process can exist in the system at the same time although they may not be all active. One processor can only execute one process at any time. If the number of processors in a system is less than the number of processes to be executed, then either a time sharing or queuing system is needed for processor sharing. The concept of processes helps to simplify the design of a complicated system. Virtually, a system can initiate a separate process to handle a particular request. A process, hence, needs only to handle a single job. Instead of considering all the job requests together in a single complicated program, system design is decomposed into two parts: the design of each individual process which handles a given job, and the management of processes. The design process is therefore greatly simplified. In DDBMS's, a transaction may involve more than one local system and hence system design is much more difficult. It is better to handle the transactions using the concept of processes.

There are three different approaches to handle transactions in a DDBMS. In Bernstein's model [BER81b], as shown in Figure 1.10, a Transaction Manager (TM) and a Data Manager (DM) reside in each site. A TM handles all the transactions entering the site and a DM handles all the queries issued by either the local or remote TM's. Obviously, the design of TM's and DM's in this approach is not easy.

In another approach, a process is created for each data access/update request entering the site either locally or remotely [BHR82] [ST080]. In this approach, creating and initializing a process to handle each remote request require too much overhead. Further, the sequence of requests within a transaction must somehow be tied together and kept in order. With respect to data access authorization, each request must be individually authenticated, which is an expensive operation.

In the third approach, used in R^* [LIN84], a process is created for a transaction and retained for subsequent use for the duration of the transaction session. There is only one setup cost for each transaction. More importantly, keeping track of the sequence of requests within a transaction is no longer a complicated job. The transactions in this research are processed in a way similar to System R^* 's approach.

1.10.6.2 Transactions and Processes in DDBLMN

Refer to the system model shown in Figure 1.9, a transaction entering the system through a local site (called the *home site* of the transaction) is processed by a TH initiated by the TM in the home site. TH may then ask the TM's in remote sites to initiate TH's to process the transaction together. All these TH's

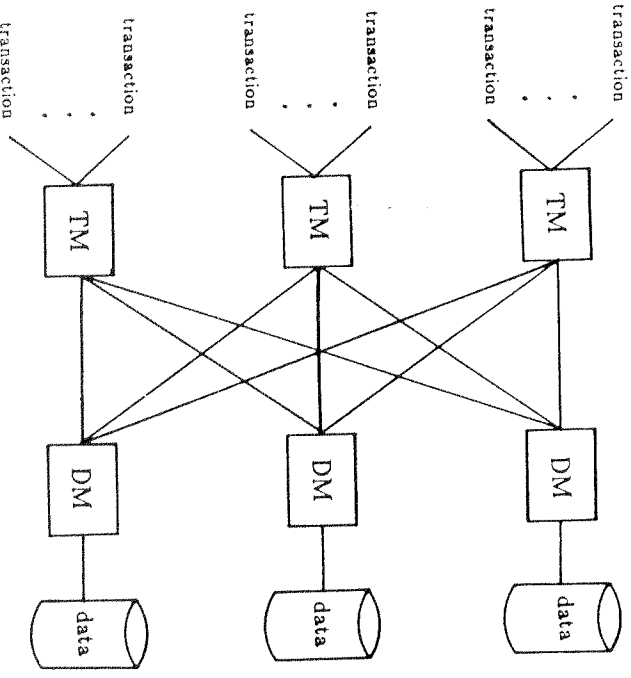


Figure 1.10 Bernstein's transaction processing model.

belong to the same transaction and are retained until the transaction is terminated. In this way, a transaction is carried out by a set of TMs cooperatively as shown in Figure 1.11.

In a DDBMS, a transaction cannot be processed independent of other transactions. In addition to the TMs, other processes are needed to coordinate the activities among the transactions. The Concurrency-Control subsystem in our model is the process in each site to handle concurrency control. Each data access or update request generated in a site is sent to the Concurrency-Control subsystem in the site, which will in turn check the consistency between the request and the current database state. If there is no consistency problem, the request will be forwarded to the local File-Server subsystem, which is another process handling the local database, and is then processed there.

1.10.6.3 Transaction Processing Strategies in DDBLMIN

A transaction consists of a sequence of independent queries. After entering the system, a transaction is first processed by the home site in a phase called the *preprocessing phase*. All necessary initializations are performed in this phase. The transaction is then broadcast to other sites for further processing. In fact, each individual query is processed in sequence by the system. The processing sequence of a transaction with three queries is shown in Figure 1.12. The preprocessing, initial-broadcast, and query processing phase will be discussed in Chapter III.

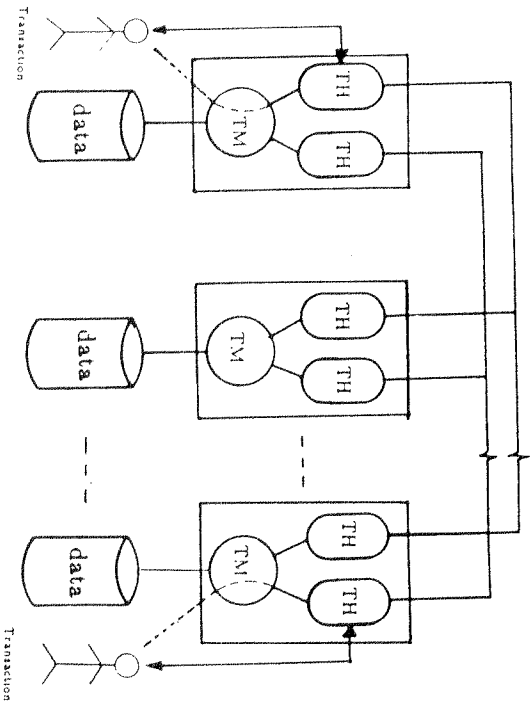
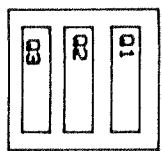


Figure 1.11 Example of two transactions being processed by the system.

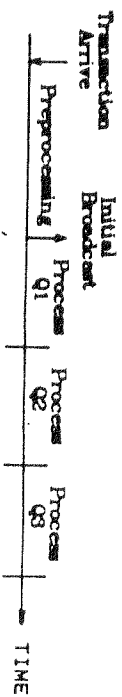
A Transaction

```

GET      R.a WHERE (R.a = S.a) AND (R.b = T.b)
UPDATE  S.a WHERE (S.a = T.a) AND (R.b = S.b)
GET      T.a WHERE (T.a = S.a)
    
```



(a)



(b)

Figure 1.12 The processing sequence of a transaction with three queries Q_1, Q_2, Q_3 (a) the transaction; (b) the processing sequence.

1.1.1 Thesis Organization

In Chapter II, the distribution design problems on DDBMIN are analyzed. A detailed cost model of file allocation problem, which is a special case of the general distribution design problem, is developed. Different cases of the problem are analyzed. Both solutions for some solvable problems and the proof for the other NP-hard problems are presented. Efficient approximation solutions for NP-hard problems are also developed.

Distributed query processing is analyzed in Chapter III, while concurrency control problem is studied in Chapter IV. A processing strategy consisting of five different phases are proposed, and various strategies for each phase are developed. A locking based concurrency control protocol for local multiaccess networks is also presented.

Under certain conditions, the query processing problem can be solved as a generalized traveling salesman problem. In Chapter V, an extensive analysis of the generalized traveling-salesman problem is presented.

In Chapter VI, simulation results to demonstrate the feasibility of the proposed DDBMS system are shown. Finally in Chapter VII, the summary of the research is given and the directions for future research are suggested.

CHAPTER II

FILE ALLOCATION

2.1 Introduction

The distribution design problem entails

- (1) the fragmentation of the database and
- (2) the allocation of these fragments with possibly redundant copies.

The objective of data distribution is to maximize the efficiency or performance of database operations subject to the system constraints. The performance can be measured by various parameters such as the overall operation cost, response time, and throughput. The constraints can be the average response time, availability, degree of redundancy, etc. The operation cost may consist of the communication cost, storage cost, and processing cost. Although the problem depends on the transaction processing strategies, solving the two problems together is very complex. In our design methodology, the distribution-design problem is solved independently from the transaction-processing problem in the database design phase.

2.1.1 Fragmentation

Fragmentation partitions a relation into several fragments, each of which has some special properties. There are two reasons that make fragmentation necessary.

(1) A relation may be too large to be stored in a single place. A large relation causes not only the storage problem but also the processing problem. For example, the buffer space in a system is limited such that there is no way to bring entire relation into the buffer to process. A considerable overhead may be induced to solve this problem.

(2) Most of the data retrievals from a certain class of users may be targeted to a small portion of the database most of the time. That is, there may exist a strong locality within a relation. A lot of the processing costs can be saved if the relations can be further partitioned such that a certain class of users may need to access a certain fragment of each relation. This may reduce processing overhead, as only the needed fragments are searched. Fragmentation can be viewed as a 'pre-selection,' which can reduce the overhead of future query operations.

As an example, a relation storing the personal information of all students in a university contains 10,000 tuples. To simplify the user interface, this information is organized as a single relation. However, most of the queries initiated from a department will only involve the portion containing information about students in the department. It is natural to decompose the relation such that each fragment only contains the information about students in a given department. Not only the processing cost is saved, but the communication cost is also saved if each fragment can be stored in the local computer of each department. An example illustrating fragmentation is shown in Figure 2.1.

There are two conditions that a fragmentation should satisfy.

- (1) **Completeness.** All data of the original relation should be mapped into fragments; i.e., a data item in the original relation must exist in at least one fragment.

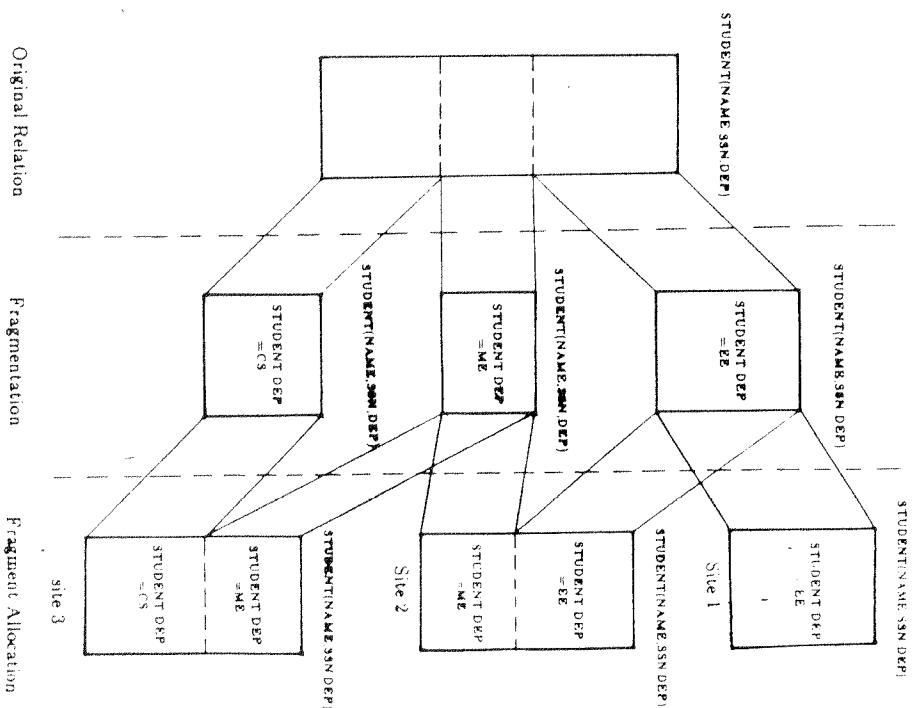


Figure 2.1 Example to illustrate the fragmentation problem

(2) **Recoverability** It should be possible to reconstruct the original relation from the fragments.

There are three basic types of fragmentations: horizontal, vertical, and mixed.

(1) **Horizontal Fragmentation** partitions the relations by tuples, in which all tuples satisfying a given condition are grouped together into a fragment. This can be done by the *selection* operation. The original relation can be recovered by the *union* operation.

(2) **Vertical Fragmentation** partitions the relations by attributes, in which a number of attributes are grouped together into fragments. Vertical fragmentation can be done by the *projection* operations. The original relation can be recovered by the *join* operations.

(3) **A mixed fragmentation** is the combination of horizontal and vertical fragmentations. A horizontal fragmentation may be followed by a vertical fragmentation or vice versa.

The Fragmentation problem in DBMS has been studied extensively by Ceri, Navathe, and Wiederhold [CER84]. We will not consider the fragmentation problem in this research. It is assumed that fragmentation is done by existing methodologies and is separated from other design problems. Consequently, it is reasonable to assume that each data fragment is a relation by itself.

2.1.2 Fragment Allocation and File Allocation Problem

The next problem after fragmentation is to allocate and place these fragments (or relations) in the system. The allocation should be based on the way that the database is queried by the users. It is impossible for a system to

collect all the access pattern down to the tuple level. We have assumed that this information is only available on the relation level. Therefore, the fragment allocation can be solved as a *file allocation problem* (FAP). The FAP entails the distribution of files (or relations) on a DBMS such that the objectives of the system are satisfied and that the overall system is operating at its most efficient level. If the constraints, requirements, query patterns, or even the database state are changing dynamically, the files must be allowed to migrate. The *file migration problem* (FMP) entails the study of the migration sequence of files to adapt to changes in the DBMS [WAL84]. Solutions developed for a general network can also be applied to a local multi access network. In this research, we will only consider the FAP as more efficient solutions can be obtained.

The FAP was originally investigated by Chu [CH69] who studied it with respect to multiple files on a multi-processor system. He modeled the problem as an integer programming problem in which the objective is to minimize the overall operating costs under the constraints of response time and storage capacity. Subsequently, a lot of work has been done in this area. There are good comparative surveys in Dowdy's [DOW82] and Wah's [WAL84] paper.

To allocate a copy of a file to a site, it has the following effects on the overall operating cost and constraints.

- (1) The communication cost of remote accesses initiated from this site is saved. It may also provide lower access cost for other nearby sites.
- (2) The storage cost, which is assumed to be proportional to the data volume, in this site is increased.
- (3) The communication cost to update this file may be increased.

- (1) The storage space in this location is reduced.
- (3) The average response time may be reduced since some remote accesses are localized.
- (6) The availability of this file is increased.

In general, the decision to allocate a specific file to a specific location cannot be made independent of the allocation of other files. The problem is formulated into an integer programming problem with the overall operating cost as the objective to be minimized. The problem is NP-hard on general networks. However, the problem is simplified on local multiaccess networks and is polynomially solvable in some special cases. There are two reasons why the problem is simplified.

- (1) The remote access/update costs on the network is site independent.
- (2) The update cost for a file is independent of the number of copies allocated if more than one copy is allocated.

In this dissertation, the different cases of the *simple file allocation problem* (SFAP), which is a special case of FAP, and the general file allocation problem with storage constraints are studied and solved. The following notation are defined.

- N - number of sites in the system;
- F - set of files in the database, $|F| = m$;
- I^f - index set of sites with a copy of file f ;
- A_j^f - query load originating at site j for file f per unit time;
- ϕ_j^f - update load originating at site j for file f per unit time;
- d - cost per unit size of query from any site to any other site;

d^f - cost per unit size of update from any site to all other sites;

α_k^f - storage cost per unit time of file f at site k ;

L^f - length of file f ;

S_k - storage capacity at site k ;

$$Y_k^f = \begin{cases} 1 & \text{if file } f \text{ exists at site } k \\ 0 & \text{otherwise} \end{cases}$$

2.2 Simple File Allocation

The SFAP considers the allocation of a single file, in which the effects of queries, updates, and data storage are represented as costs [RAM83]. It has been shown to be NP-hard [ESW74] [GAR79], which means that the computation time for all known optimal algorithms for this class of problems increases at least exponentially with the problem size. Numerous exhaustive optimal algorithms and suboptimal heuristics have been studied [WAh84a] [DOW82].

In this section, a solution to the SFAP for local multiaccess networks is proposed. Since only one broadcast is required to either lock, unlock, or update a file, the problem is found to be polynomially solvable with respect to the number of sites.

In the following cost function for a given placement I , the index f is dropped because the problem is defined with respect to a single file.

$$C(I) = \sum_{j=1}^N \lambda_j d(1 - Y_j) + \sum_{j=1}^N \phi_j^f d^f \left(1 - Y_j \prod_{k=1}^N (1 - Y_k) \right) + \sum_{k=1}^N \alpha_k^f Y_k \quad (2.1)$$

The first and last terms on the RHS of the above equation represent the query and storage costs, respectively. The second term accounts for the update costs.

When multiple copies of a file exist in the system, each unit of update incurs a constant cost of d' . However, when only one copy of the file exists, updates originating from the site with the copy do not incur any cost because they do not have to be broadcast.

This problem can be formulated as an integer program with Eq. (2.1) as the cost function and constraints that at least one copy exists and Y_k is either 0 or 1.

Three different cases are considered in this section: the conventional SFAP, SFAP with availability constraints, and SFAP with average delay-time constraints.

2.2.1 Conventional SFAP

The solution to this problem can be simplified by considering the single-copy and the multi-copy cases separately.

For the single-copy case, the cost function becomes

$$\begin{aligned} C(I) &= \sum_{j=1}^N [(\lambda_j d + \phi_j d')(1 - Y_j) + \sigma_j Y_j] \\ &= \sum_{j=1}^N (\lambda_j d + \phi_j d') + \sum_{j=1}^N (\sigma_j - \lambda_j d - \phi_j d') Y_j \end{aligned} \quad (2.2)$$

In addition to the fixed cost in Eq. (2.2), the cost of not placing the file at site j is $\lambda_j d + \phi_j d'$, while the cost of placing the file there is σ_j . Thus the file should be placed at site j if the cost difference $(\sigma_j - \lambda_j d - \phi_j d')$ on this site is the minimum. The optimal allocation can, thus, be found in $O(N)$ time.

For the multi-copy case, at least two copies must be allocated. The cost function is

$$C(I) = \sum_{j=1}^N [\lambda_j d(1 - Y_j) + \phi_j d'] + \sum_{k=1}^N \sigma_k Y_k \quad (2.3)$$

$$= \sum_{j=1}^N (\lambda_j d + \phi_j d') + \sum_{j=1}^N (\sigma_j - \lambda_j d) Y_j$$

The first term in the above equation is constant. The optimal allocation can be found by allocating a copy at site j if $(\sigma_j - \lambda_j d)$ is negative. This will decrease the overall cost, $C(I)$. In case that none or one of the cost differences is negative, two copies with the minimum cost differences are selected. The optimal allocation can hence be found in $O(N)$ time.

The global optimum is obtained by comparing the costs under the single-copy and multi-copy cases. Therefore, the SFAP is optimally solvable in polynomial time. As an example, consider a system with four nodes, query rates $\lambda_j = [6, 7, 4, 5]$, update rates $\phi_j = [3, 6, 2, 4]$, and storage costs $\sigma_j = [3, 5, 2, 5]$. The per unit query and update costs, d and d' , are assumed to be unity. In the single-copy case, $[\sigma_j - \lambda_j d - \phi_j d'] = [-6, -8, -4, -4]$. The single copy should be allocated at site 2 with cost 29. In the multi-copy case, $[\sigma_j - \lambda_j d] = [3, -2, -2, 0]$. Copies should be allocated to sites 1, 2, and 3 with a total cost 30. Comparing the single-copy and multi-copy cases, a single copy should be placed at site 2.

2.2.2 SFAP With Availability Constraints

The SFAP can also be solved efficiently when an availability constraint is included. As mentioned in Chapter 1, the availability is defined as the probability that at least one copy of a file is available. Since the network is reliable and the failure rates of sites are assumed to be identical, the availability constraint can be expressed into a minimum number of copies of the file. The locations for these copies can be determined as above.

There is no difference between the the SFAP with an availability constraint and the conventional SFAP if the minimum number of copies of the file is one, and the number of copies should be more than two, then only the multi-copy case needs to be considered.

A copy of the file is allocated at each site with negative (σ_j, λ_j, d_j). Additional copies are allocated such that the availability constraint is met and the cost increase is minimum. The complexity of this algorithm is also $O(N)$.

2.2.3 SFAP With Average-delay Constraints

The SFAP becomes NP-hard when an average-delay constraint is imposed.

Given that the average delay of messages on the multiaccess bus is a constant, c , the average delay for all queries and updates must satisfy:

$$\sum_{j=1}^N \left[\lambda_j c(1 - Y_j) + \phi_j c \right] \leq D' \quad \text{or} \quad \sum_{j=1}^N \lambda_j c(1 - Y_j) \leq D \quad (2.4)$$

In Eq. (2.4), updates are assumed to be always broadcast. The SFAP with the above average delay constraint can also be separated into the single-copy and multi-copy cases. The single-copy case is solvable in $O(N)$ time. For the multi-copy case, the cost formula in Eq. (2.3) can be rewritten as

$$C(I) = \sum_{j=1}^N (\phi_j d' + \sigma_j) + \sum_{j=1}^N (\lambda_j d - \sigma_j)(1 - Y_j).$$

Since the first term on the RHS is constant, the optimization problem for the multi-copy case can be written as

$$\text{maximize } C(I) = \sum_{j=1}^N (\sigma_j - \lambda_j d)(1 - Y_j). \quad (2.5)$$

$$\text{subject to } \sum_{j=1}^N \lambda_j (1 - Y_j) \leq \frac{D}{c} \quad \text{and}$$

$$1 - Y_j = 0 \quad \text{or } 1$$

This problem is reducible from the 0-1 knapsack problem with N objects by the variable transformation $Z_j = 1 - Y_j$, where the profit and weight of the j 'th object are $(\sigma_j - \lambda_j d)$ and λ_j , respectively. However, the profits may be negative for some of the objects, and they are excluded from consideration in the optimization. The problem can be solved by dynamic programming algorithms or fully polynomial-time approximation schemes.

2.3 General FAP With Storage-Capacity Constraints

The general FAP considers the allocation of multiple files under such design requirements as delay, storage capacity, parallelism, and availability. As a simplification, each query is assumed to access a single file, and all file accesses are assumed to be independent. Only the storage capacities of sites are considered in the following formulation.

The problem with storage-capacity constraints can be formulated as:

$$\begin{aligned} \text{minimize } C(I^1, \dots, I^F) &= \sum_{j=1}^F \lambda_j^j d(1 - Y_j^j) + \\ & \sum_{j=1}^F \phi_j^j d' \left(1 - Y_j^j \prod_{k \neq j} (1 - Y_k^k) \right) + \sum_{j=1}^F \sigma_j^j Y_j^j \end{aligned} \quad (2.6)$$

subject to:

- Given a finite set U of m objects and a knapsack of size B , there are a size $s(u)$, $s \in Z$ and a profit $p(u)$, $p \in Z'$ for each $u \in U$, where Z is the set of positive real numbers. The knapsack problem searches for a subset $U' \subseteq U$ that maximizes $\sum_{u \in U'} p(u)$ such that $\sum_{u \in U'} s(u) \leq B$.

$$(a) Y_j^f = 0 \text{ or } 1, \quad j = 1, \dots, N, \quad f \in F$$

(each site contains at most one copy of each file);

$$(b) \sum_{j=1}^N Y_j^f \geq 1, \quad f \in F$$

(at least one copy of each file exists in the system);

$$(c) \sum_{f \in F} L^f Y_j^f \leq S_j, \quad j = 1, \dots, N$$

(the storage capacity at each site is not exceeded).

The above problem is non-linear, but can be linearized easily by using different index variables [GEO72]. This problem is NP-hard as shown in the following theorem.

Theorem 2.1: The optimization problem defined by Eq. (2.6) is NP-hard.

Proof: We show that the 0-1 knapsack problem reduces to this problem in polynomial time. Given an instance of the 0-1 knapsack problem, an instance of the problem in Eq. (2.6) can be formed with the following parameters: $N = 2$; $F = U$; $L^f = s(f)$ for $f \in F$; $S_1 = \sum_{f \in F} L^f$; $S_2 = B$; $\sigma_1^f = 0$ for $f \in F$; $\lambda_1^f = a$ large

constant for $f \in F$; ϕ_1^f, λ_2^f and σ_2^f are chosen such that $\lambda_2^f d - \phi_1^f d' - \sigma_2^f = p(f)$ for $f \in F$. Basically, the first site has large enough capacity to hold a copy of all files in F . The query and storage costs are chosen such that a copy of all files are allocated at site 1. Therefore, the problem becomes the packing of the second site, which is the standard knapsack problem. \square

2.3.1 Optimal Solutions of FAP With Storage Constraints

Although the problem in Eq. (2.6) resembles the packing of $m \cdot n$ copies into m knapsacks, there is a subtle difference here. The profit of allocating the first copy of file f to site j is $(\lambda_j^f d + \phi_j^f d' - \sigma_j^f)$, while the profit of allocating the second copy of file f to site $k, k \neq j$, is $(\lambda_k^f d - \phi_j^f d' - \sigma_k^f)$. The extra term $\phi_j^f d'$ in the profit of the second copy offsets the additional profit incurred when the first copy is allocated. The profit of allocating other copies of file f to site $q, q \neq j$ and $q \neq k$, is $(\lambda_q^f d - \sigma_q^f)$. Note that these profits may be negative. Due to the above and since one or more copies must exist in the system, the problem is decomposable into multiple independent knapsack problems only if at least two copies have been allocated in the system. The overall profit can then be maximized by independently optimizing the profit of each site with the remaining capacity. An optimal algorithm, therefore, consists of enumerating the allocations of the first two copies of each file, and solving N knapsack problems for each combination. The structure of the state-space tree is as follows.

- (1) The root is in level 0.
- (2) In the first m levels, the first copy of all files are allocated. Level $f, 1 \leq f \leq m$, represents the allocation of file f . Each file can be allocated to one of the N sites. Thus the degree of the tree in the first m levels is N . If a file cannot fit in a site, the subtree representing this allocation is considered infeasible and terminated.
- (3) The allocation of the second copy of each file is carried out in the $(m+1)$ st to the $2m$ th levels of the tree. Since the second copy must not be allocated to the site containing the first copy and may not have to be allocated in the DDBMS, the degree of branching in each level is N . If a copy cannot fit in a

site, then the subtree representing this allocation is considered infeasible and terminated. Further, if an allocation of the second copy at site k results in negative value of $(\lambda_k^j d - \sigma_k^j)$, then the second copy should not be allocated at site k , and the corresponding subtree is also terminated.

- (4) For all active nodes left in the state-space tree in level $2m+1$, the allocation of the remaining copies is solved as N single-knapsack problems. For a particular site, the files to be considered for packing into the remaining capacity are those that have not been allocated at this site and those with the second copy allocated in the system. The algorithm for solving the conventional 0-1 knapsack problem has to be modified to take into account the negative profits.

Bounding criteria similar to those used in branch-and-bound algorithms can be developed here [GEO72]. In particular, an upper bound on the profit can be computed for each node in the state-space tree by a linear program or a greedy algorithm without the integrality constraints. The maximum profit of the currently available feasible solutions is kept in the incumbent. If the upper bound of a node is smaller than the incumbent, then the subtree originating from this node is pruned because it cannot possibly lead to an optimal solution. An example illustrating the allocation of two files in a system with two sites using a depth-first search is shown in Figure 2.2. In computing the upper bound for a node in the branch-and-bound tree, each site is treated as an independent knapsack. Given that $i, i=0, 1, 2$, copies of a file to be included in the knapsack have been allocated in ascendant levels of this node, the profit of the file to be used is that of the $(i+1)$ th copy. A greedy algorithm computes the upper bound by assigning the files to each knapsack in descending ratios of profit to length until the knapsack is full. To keep the knapsack full, the last file may be

Cost of query, update, and storage				
	file 1	file 2	file 1	file 2
$\lambda_j^i d$	8	9	6	7
σ_k^j	1	1	1	1
c_k	4	6	4	6

Storage limit of sites	
site 1	site 2
5	6

Length of files	
file 1	file 2
3	2

Profit of first and second copy				
	site 1	site 2	site 1	site 2
First copy	3	4	3	2
Second copy	3	2	1	0

Allocations:

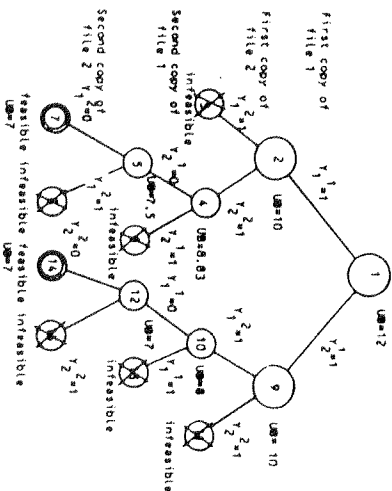


Figure 2.2

Example to illustrate the general file allocation problem

partially allocated. The worst-case number of knapsack problems to be solved in this approach is N^{2m} , which is much better than an exhaustive enumeration.

2.3.2 Heuristic Solutions of FAP With Storage Constraints

When the problem size is large, an optimal solution may not be possible to get in a reasonable amount of time. More efficient approximate solutions are needed.

The algorithm in Section 2.3.1 can be easily modified to get good approximate solutions. The knapsack problems in the algorithm can be solved heuristically, instead of optimally. Although the search method in Section 2.3.1 requires a total of N^{2m} knapsack problems to be solved, it is still a useful algorithm because the FAP can be solved in design time instead of in real time.

2.4 Summary

In summary, FAP on local multiaccess networks is simplified by the broadcast capability and the property that the communication costs are site independent. The solution algorithms for the SFAP are efficient and can be evaluated in either polynomial or pseudo-polynomial time. These algorithms can be applied in real time when the access rates change. However, optimal algorithms for the general FAP are complex and should be used in the design stage only. The algorithm can be easily modified into a less complex heuristic algorithm for large problems.

CHAPTER III

DISTRIBUTED QUERY PROCESSING

As mentioned in Section 1.10, a transaction is processed in a strategy that consists of preprocessing phase, initial broadcast, and query processing phase. The preprocessing phase and initial broadcast are discussed in Section 3.1. The query processing phase is discussed in the rests of this chapter.

3.1 Preprocessing Phase and Initial Broadcast

Preprocessing Phase

The operations which need to be performed in this phase are:

- (1) **Administration.** These include authorization, accounting, etc.
- (2) **Query decomposition and transformation.** Users of a database usually query the database using a high level 'user friendly' query language. Direct execution of queries presented in the original form is usually inefficient. A

decomposition and a transformation of the given query may result in more efficient execution. For example, in target-data identification, a sequence of basic operations are applied to the database to identify the needed data.

The size of the intermediate results may be different in different processing sequences. It is generally more efficient to apply unary operations, such as selections and projections, as early as possible.

(3) *Related-site identification and materialization.* *Related sites* of a transaction are the sites participating in the transaction processing. Materialization identifies the specific copies of relations needed by a transaction. A relation in a site is *materialized* if it is identified in the materialization. The related sites of a transaction in DDBLMN are the sites where at least one materialized relation is stored.

Initial Broadcast

After preprocessing, the transformed transaction with the necessary information is broadcast to all related sites. This is called *initial broadcast*. Each related site, then, initiates a TH to participate in the transaction processing.

Since the network is using a contention based bus access protocol, it is remotely possible that a transaction can never perform its initial broadcast. Further, a transaction is not known to other related sites until its initial broadcast has been carried out. Hence if the network can support priority messages, initial broadcasts should be given higher priorities for the transaction.

3.2 Introduction to Distributed Query Processing Problem

Minimizing processing cost is the most important goal of query processing.

A large amount of communication traffic may be involved in distributed query processing (DQP). The communication cost is usually proportional to the volume of transmitted data in most DDBMS's, including DDBLMN. Hence the objective of query processing in these systems is to reduce the total data volume to be transmitted.

The main difficulty of query processing is the target data identification, which involves the identification of the set of physical fragments (relations) referenced by the query (*materialization*), the selection of the order of execution of operations, and the selection of the method for executing each operation.

The *raw-target* in the execution of a target-data identification is defined as the smallest portion of the database that contains the 'target-data'. The procedure of target-data identification is like applying a sequence of 'cuts' on the 'raw-target' to get the target-data. Initially, the raw-target is the entire database. After applying an operation, the portion that is known not to contain the target-data is cut off. After a sequence of cuts, the target-data is finally identified. This procedure in DDBMS is known as *reduction* and is shown in Figure 3.1. One difference between the procedure in the example and the real DQP is that the size of a raw-target may be increased after a join operation. Obviously, the query processing cost is a function of the total volume of raw-targets in all the intermediate steps. The selection of the order of reductions becomes the major task of DQP.

3.2.1 Redundant Materialization

In most of the previous studies, the materialization is assumed solved independently and is non-redundant, in which one copy of each relation is identified in processing the query. A possible alternative is the redundant materialization, in which more than one copy of each relation may be identified. The redundant materialization has the following advantages over the non-redundant materialization.

(1) Non-redundant materialization is a special case of redundant materialization. For each schedule obtained from a system using non-

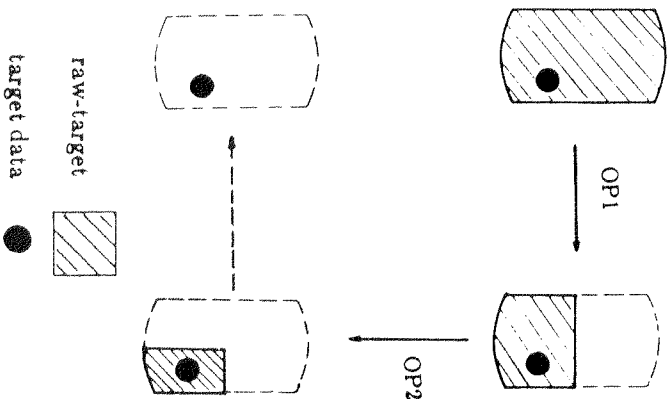


Figure 3.1 An Example to illustrate the procedure of target data identification.

redundant materialization, a schedule with less or equal cost can always be found from a system using redundant materialization.

- (2) Redundant materialization is more flexible. In a DBMS, a local database may have to answer queries initiated from more than one transaction at the same time. A transaction has to wait if the local database is busy. With redundant materialization, the availability of data is higher.

- (3) Redundant materialization may have less query processing cost in a heuristic query processing. Whether a piece of the raw-target contains the target-data or not usually depends on the relationship among different components of the raw-target. For example, in the following query

GET R.a WHERE (R.a = S.a) AND (R.a = T.a),

the target-data is a portion of R.a and the initial raw-target consists of three components R, S, and T. They can be easily reduced into R.a, S.a, and T.a in the beginning. To perform further reduction, the relationship between R.a and S.a, and the relationship between R.a and T.a should be known. If relations R, S, T are stored at different sites, then there is no way to know the relationship among them to perform any reduction in each local site. Without an exchange of information among sites, the smallest raw-target which can be identified is the set containing R.a, S.a, and T.a. If some of them are stored at the same site, say R and S, then the raw-target can be reduced to a set containing T.a and the portion of R.a whose values are also found in S.a. From the above example we can see that the more relations a site has, the more relationships among target-data components for data reduction it can have. Using redundant materialization in target-data identification, the materialized relations in each site are generally more than that of using non-redundant materialization, and consequently, more

relationships among the raw-target components are available for the raw-target reduction. Hence, the volume of final raw-target in a site can be smaller. For example, in the same query as above, the volume of raw-target in a site storing R, S, and T should be the smallest (even the target-data can be identified in this case) among all sites.

Although the redundant materialization may have smaller communication cost, more local processing cost may be involved. A system emphasizing on local processing cost may not be adequate for redundant materialization.

3.2.2 Processing of Joins

Among three most frequently used operations (selection, projection, and join), the join operation is the only one that is a binary operation and that may expand the volume of raw-targets. The efficient processing of join operations becomes the most important job of target-data identification in both CDBMS and DDBMS.

3.2.2.1 Multi-site Joins and Semi-joins

Processing a join operation is more complicated in DDBMS since the two operands of a join operation may reside in two different sites. One approach to process such a join operation is to send one relation in one site to the other site and execute the join operation in the second site. The cost of this approach is a function of the size of entire relation sent. A better approach proposed by Bernstein and Chiu [BER81a] [BER81c] is *semi-join*. They observed that when two distributed relations are joined, one of the operand relations can be reduced in size by deleting those tuples that do not appear in the final result. This may be achieved by sending values of the joining attribute of the second relation to

the first one and performing a join there. This reduction effect is called a *join restriction* [YAO79]. The join is completed by sending the reduced form of the first relation to the second one. Referring to the raw-target model, a semi-join is essentially cutting the raw-target at one side first, then sending the remaining raw-target to the other side. As an example, relations $R(a, b)$ and $S(b, c)$ in Figure 3.2(a) are stored at sites A and B, respectively, and are to be joined on the joining attribute, b. In the first approach, as shown in Figure 3.2(b), relation S is sent to site A and is joined with relation R there. The communication cost is a function of the size of relation S. The semi-join approach is shown in 3.2(c), in which the unique values of attribute R.b is first projected and sent to site B, then relation S is reduced by eliminating those tuples whose values on attribute b are not in attribute R.b. Finally, the remaining tuples of relation S are sent to site A and are joined to R there. The communication cost is a function of R.b and a portion of S.

When multiple relations are to be joined, the most promising approach currently is to first execute a sequence of semi-joins among relations to reduce the size of relations, then the remaining fragments are transferred to a single site at which the complete join is performed. Since the major overhead is due to semi-joins, the minimization of this overhead is the objective in most DQP studies. A comprehensive approach of query processing that incorporates semi-joins has been developed by Apers, Hever, and Yao [HEV79a] [APE83].

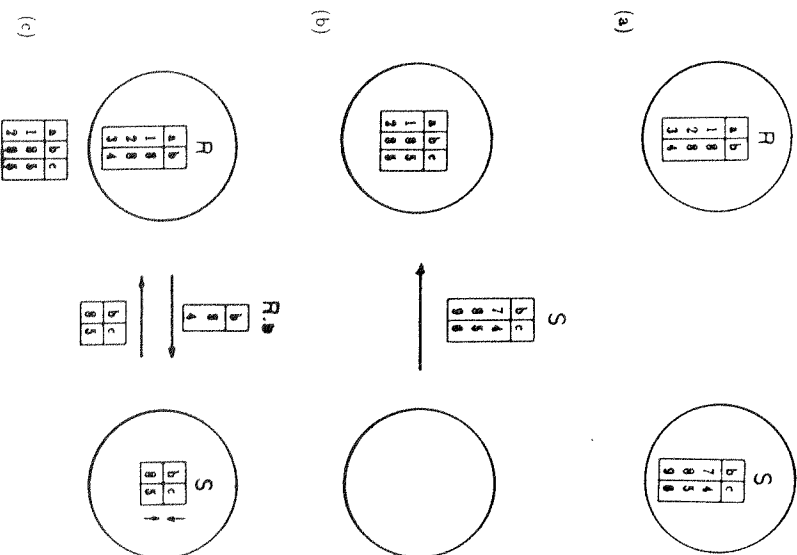


Figure 3.2 An example to illustrate the procedure of join and semi-join
 (a) relation R and S; (b) sending relation S to join with R;
 (c) semi-join.

3.2.2.2 Size Estimation and Selectivity

It's obvious that the performance of a DQP algorithm highly depends on the estimation of the sizes of raw-targets. The choice of a reasonable estimation algorithm is, therefore, extremely important.

Suppose that relations R_1 and R_2 are single-attribute relations and the values of the common attribute 'a' are uniformly and independently distributed on the relations. The size of R_2 after joined with R_1 has to be estimated.

By letting $p_{i,a}$ be the probability that a value appears in R_i , $i = 1, 2$, then $p_{i,a}$ is called the *selectivity* of R_i on attribute a. Since the expected total number of distinct values existing in attribute a of both relations is $|a| \times p_{1,a} \times p_{2,a}$, where $|a|$ is the cardinality of the domain a. The size of the reduced R_2 can be estimated to be $|a| \times p_{1,a} \times p_{2,a} \times w_a$ or $|R_2| \times p_{1,a} \times w_a$, where $|R_2|$ is the cardinality of R_2 , and w_a is the width of attribute a. Essentially, the size of R_2 is reduced by a factor of $p_{1,a}$.

In case of joining multiple-attribute relations, not only the size of the joined attribute and the size of the joined relation are reduced, but the size of other attributes are also reduced. This is called an *indirect semi-join*. In the above example, suppose that R_2 is a relation with two attributes, a and b. After semi-joining R_1 with R_2 , the cardinality of R_2 can be estimated as $|R_2| \times p_{1,a}$. The cardinality of R_2 was estimated by Bernstein and Chin [BERNSTA], which can be demonstrated in the following *ball-color* problem. In this problem, there are n balls with m different colors. It is necessary to find the expected number of colors if t balls are randomly selected from the n balls. The correspondences between the ball-color problem and semi-joins are as follows. The value n corresponds to the number of tuples of R_2 before the semi-join, m is the number in distinct values of R_1 projected on attribute b before the semi-join, and t

corresponds to the number of tuples of R_i after the semi-join. The expected number of colors of the l selected balls is

$$m \times \left[1 - \prod_{i=1}^l \left(\frac{n((m-1)/m) - i + 1}{n - i + 1} \right) \right].$$

The computation of this formula is expensive. An approximation is needed and was given by Bernstein and Chiu. This solution is the same as the solution in Yao's paper [YAO77] for estimating the hit ratio of the block-access problem.

3.2.2.3 Static Semi-join Scheduling

Yao, Heyner, and Yperis [APE83] [HEY79a] proposed a *static semi-join scheduling* strategy to solve the DQP. In which, the processing costs of different semi-join schedules are estimated first based on their selectivities; the schedule with the least estimated cost is selected; and the query is executed according to the selected schedule. The problem is NP-hard even for simple queries with one common joining attribute [HEY79b]. The main advantage of this approach is that the optimal schedule of a certain type of queries can be precompiled and reused later. There is no need to produce a schedule for each query processed in real time. On the other hand, there are some disadvantages associated with the static semi-join scheduling.

- (1) An accurate estimation is hard to get. The computation of selectivities does not consider the dependency among attributes. In fact, attributes in a real database are usually dependent, although the effects of semi-joins on dependent attributes are difficult to estimate. Further, the assumption on uniform distribution may not be realistic. Lastly, the indirect semi-join effect complicates the problem. In most current research on static semi-join

scheduling, the indirect semi-join effect is simply not taken into account.

- (2) The errors in estimation usually propagate and accumulate after a sequence of semi-joins. This is especially true when the model of estimation is simple.
- (3) A static schedule is inflexible. Once a semi-join schedule is determined, the schedule is generally not allowed to change during the execution. The entire transaction is blocked if the needed relation in any stage is not available. The probability of being blocked is increased as the number of required relations increases or the system load increases.
- (4) The computational overhead of static strategies is high. Although heuristic algorithms are generally used to determine semi-join schedules, the computational time is still very high, especially when more accurate solutions are needed. This overhead can only be reduced by using precompiled semi-join schedules in a relatively stable system.

3.2.2.4 Dynamic Semi-join Scheduling

An alternative is to use *dynamic semi-join scheduling* [YE83] [SAC84] [WAH84b], in which each individual step in the schedule is determined immediately before the execution of that step based on the current state of the database. Such a dynamic scheduling algorithm is a greedy heuristic algorithm. Some of the problems associated with static scheduling do not exist in dynamic scheduling.

- (1) There is no accumulated estimation error. In each step, at most one estimation is made in dynamic scheduling.

(2) Dynamic scheduling is more flexible, and the probability of being blocked is lower. Whenever the required relation is not available during the execution, the schedule can be changed. Another way is to only schedule semi-joins for available relations.

The success of dynamic semi-join scheduling relies on the selection of the joining attribute. In greedy heuristic, a function called *heuristic function* is designed to measure the possibility of leading to the optimal solution for each candidate. Therefore, the key problems are the design of the heuristic function and the collection of database statistics. Designing a good heuristic function is as complicated as designing a good static scheduling algorithm. There is no easy way to collect statistics on general networks. Status exchange on a distributed system not only involves a lot of communication overhead, but also slows down the query processing. Further, real-time statistics exchange is almost impossible in general networks because the system statistics may be already outdated when it is collected. Therefore, a necessary condition for a DDBMS to employ a dynamic semi-join scheduling algorithm is that an efficient system statistics collection mechanism must be provided. This difficulty has prohibited the use of dynamic schedules in general networks. Fortunately, the broadcast capability of local multiaccess/broadcast bus allows an efficient distribution of status information. Hence, dynamic semi-join scheduling is employed in DDBLMN.

3.3 Previous Work

Most of the previous studies on distributed query processing were based on relational databases in a non-broadcast system. A materialization is assumed for each query, the order of processing is optimized, and the operations are clustered into local actions. Further, the transmission cost per unit data is

assumed to be constant for any two points in the network, and the cost of local processing is negligible. Since the DQP is NP-hard, heuristic algorithms were generally used [WON77]. The semi-join approach was first incorporated into DQP strategies in SDD-1 system by Bernstein and Chin and was shown to be more efficient than full joins [BER81a] [BER81c]. Yao et al. improved this approach with a static semi-join scheduling strategy. Considerable amount of research have been carried out on finding good static semi-join scheduling strategies.

Further, studies have been made to improve semi-joins. Yu et al. reduced the transmission cost by incorporating complement transmissions into semi-joins [YT82]. The *complement* of an attribute is defined as the set of values absent from a given attribute. In semi-join operations, the complement of an attribute is sent if its cardinality is smaller than that of the original attribute. Nonetheless, there is no efficient algorithm to determine the optimal semi-join sequence, even for simple queries. Data compression techniques have also been applied in semi-joins. For instance, a bit vector indicating the absence or presence of a value can be sent instead of the attribute if the size of the bit-vector is smaller than the size of the attribute [GOU81]. Suboptimal techniques for semi-joins with smaller optimization overhead have also been studied [KR184]. Chen and Li employed a graphical approach to identify and eliminate redundant transmissions in a semi-join schedule [CHE84]. Query processing on a star network was studied in [KER82]. Static query processing strategies for local area broadcast networks and address ring were developed in Hevner, Wu, and Yao's paper [HEV85]. A good survey was given by Yu and Chang [YU84b].

Chu and Hurley have developed a unified approach that considered both the local processing and transmission costs [CHU82]. A query-tree model for

selecting the transmission sequence and the sites for executing a set of subqueries was proposed.

In short, previous studies on DQP assumed that attributes were independent, ignored the effect of semi-joins on non-joining attributes, and used a non-redundant materialization as well as a static strategy.

3.4 Query Processing Strategy in DDBLMIN

In this section, we describe a solution of the DQP with possibly redundant materialization and a dynamic strategy. An attribute is first selected for broadcasting. Semi-joins are performed and the statistics of the resulting database are collected at all sites. Based upon these statistics, the sites cooperatively select the next attribute to broadcast. The process is repeated. This strategy has been proposed before [YU83] [SAC84], but was hampered by the high cost of collecting statistics after each semi-join. The use of a broadcast network allows the collection of statistics at a reasonable cost. The materialization can be either non-redundant or redundant. The system designer should make the decision based on the ratio of the communication cost and the actual local processing cost if it is important. Since non-redundant materialization is a special case of redundant materialization, only redundant materialization is discussed here. We assume that an independent materialization algorithm is provided if the non-redundant materialization is used.

3.4.1 Five Phases in Query Processing

A query processing strategy consisting of five phases is executed at each related site: Concurrency-Control (CC), Local-Processing (LP), Global Semi-join (GSJ), Relations-Transmission (RT), and Post-Processing (PP) phases.

3.4.1.1 Concurrency-Control (CC) Phase

Our proposed concurrency control algorithm utilizes locks. Hence, the relations used in a query must first be locked before any processing is carried out. Locking can be achieved in a distributed fashion by utilizing the transaction information broadcast in the Initial Broadcast Phase. If all relations are to be locked initially before the transaction is processed, this phase is completed at each related site when all needed local relations are not locked by other transactions. However, if the precedence relationship is resolved whenever conflicts are detected, as proposed in the concurrency control algorithm in Chapter IV, then this phase has to be repeated before each semi-join broadcast.

3.4.1.2 Local-Processing (LP) Phase

The LP phase may be started immediately at a related site as soon as the CC Phase is completed. Selections, projections, and local joins are performed on all materialized relations, which are all the relations referenced to by the query, in this phase. After local processing, the relations needed by the query in a related site are joined together into a single relation called the *site-relation*. An attribute in a site-relation is called a *site-attribute*. For clarity, the original relations and attributes are referred to as the *original-relations* and *original-attributes*. The site-relations and site-attributes are used in the discussion of the query processing strategy. The use of the original-attributes will be discussed after the various phases are presented.

Depending on the lock requests received from the bus, the results of the LP phase for a transaction may have to be retracted if a conflicting lock request is received before the first semi-join broadcast for this transaction is sent. This problem can be avoided by broadcasting a message to inform all sites that the LP Phase for a transaction has started, so other sites will not send conflicting lock requests. The synchronization of the LP Phase will be discussed in the next section. The LP Phase for all related sites should be completed before the next phase can begin. This requires the identification of the slowest site and will be discussed in the next section.

3.4.1.3 Global Semi-join (GSJ) Phase

In this phase, semi-joins on site-attributes are carried out. An initial site-attribute is broadcast, and site-relations at other related sites are semi-joined with the broadcast site-attribute. The statistics of the site-attributes are then collected at each related site. Based on measures to be discussed in Section 3.4.4, the site-attribute with the minimum *heuristic value*, called the *minimum attribute*, is selected as the next site-attribute to be broadcast. An algorithm to identify the minimum site-attribute is discussed in the next section. A semi-join may be blocked because one or more relations containing the attribute to be semi-joined are locked by other transactions. To avoid the delay of waiting for these relations to be unlocked, the processing order of semi-joins may be rearranged, so attributes that are not locked can be processed first. The process of selecting and broadcasting minimum attributes is repeated until all site-attributes are broadcast. Two techniques, complement broadcast and non-profitable attribute identification, can be used to reduce the communication

cost and are built into the scheduling algorithm.

Complement Broadcast

Semi-joins can be improved by the use of complement attributes [Y182]. The *initial complement* denotes the complement of a site-attribute before the GSJ Phase. As global semi-joins are carried out, the size of a site-attribute will decrease, while the size of its complement will increase. The size of the complement of an updated site-attribute will always be larger than that of its initial complement. Figure 3.3 shows the relationships among the complement, the initial complement, and the site-attribute itself. When the cardinality of the updated site-attribute is smaller than that of its initial complement, a *normal broadcast* that broadcasts the updated site-attribute should be used; otherwise, a *complement broadcast* that broadcasts the initial complement should be used. A semi-join can be carried out with either the site-attribute or its complement.

Non-Profitable Attribute Identification

Since the materialization is redundant, the information contained in one site-attribute may be a subset of the information contained in another site-attribute at a different site. For example, if R_1 is stored in sites 1 and 2, then the information of the site-attributes at sites 1 and 2 will overlap to each other. More than one copy of R_1 may be joined together. The result of joining redundant relations is equivalent to that of joining single copy relations. Because the result of equi-joining a relation to itself is equivalent to the relation itself, and join operations are commutable. However, redundant information had better not to be broadcast in GSJ phase. The following definitions are defined for the convenience of identifying the non-profitable attributes in GSJ phase.

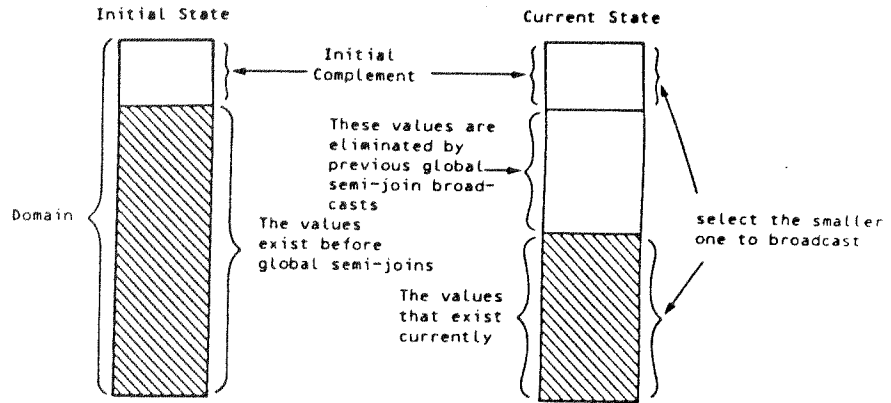


Figure 3.3 The set of values of an attribute, its complement, and the initial complement.

The *component set* of site-attribute j in site-relation R_x is defined as the set of original-relations in site x that contain site-attribute j . For example, let R_1 with attributes a and b ($R_1.a, R_1.b$) and R_2 with attributes b and c ($R_2.b, R_2.c$) be stored in site x . After the LP Phase, a site-relation (R_x) with site-attributes $R_x.a, R_x.b$, and $R_x.c$ is created by joining R_1 and R_2 . The component sets of $R_x.a, R_x.b$, and $R_x.c$ are $\{R_1\}$, $\{R_1, R_2\}$ and $\{R_2\}$ respectively.

The *broadcast component set of domain j* is defined as the union of the component sets of site-attributes in domain j that have been broadcast. In GSJ phase, a site-attribute j does not have to be broadcast if its component set is a subset of the broadcast component set of domain j , because the information contained in this site-attribute is already known to all sites. After a site-attribute has been broadcast, the corresponding broadcast component set is updated at all sites. For example, suppose R_1 is stored in sites x and y , R_2 is stored in site y , and R_3 is stored in site z . Assuming that these relations have one common site-attribute j , then site-attribute $R_x.j$, with component set $\{R_1\}$, is constructed at x ; site-attribute $R_y.j$, with component set $\{R_1, R_2\}$, is constructed at y ; and site-attribute $R_z.j$, with component set $\{R_3\}$, is constructed at z . If $R_y.j$ is broadcast first in the GSJ Phase, then the broadcast component set of domain j is updated from empty to $\{R_1, R_2\}$. Since the component set of $R_x.j$ is one of its subset, $R_x.j$ does not have any reduction effect on $R_x.j$ and does not have to be broadcast.

3.4.1.4 Relations-Transmission (RT) Phase

In this phase, the resulting fragments of the site-relations are broadcast sequentially to the *post processing sites*, where the complete join will be performed. The problem here is the determination of the order of broadcast. Since semi-joins may not execute a join completely [BER81a], an appropriate broadcast order may result in some reduction effects (join restrictions) on the unsent fragments. However, these effects are insignificant, especially when the number of semi-joins in the GSJ Phase is large. Therefore, the site-relations may be broadcast in an arbitrary order. Of course, site-relations at the post processing sites do not have to be broadcast.

3.4.1.5 Post-Processing (PP) Phase

In this phase, the full join and subsequent operations are executed, and the results are sent to the site from which the query originates. The process of determining the post processing sites is different for retrievals and updates. For a retrieval, the set of post processing sites must be selected such that the total overheads of transmission in the RT Phase and transmission of results from the post processing sites to the originating site is minimum. Since the set of post processing sites must be determined before the RT Phase begins, and the difference in overheads of processing at different post processing sites is usually small, the originating site can be chosen as the post processing site arbitrarily. For updates, the results of the post processing sites must also be sent to all sites containing relations to be updated. By an argument similar to that of retrievals, the originating site is chosen as the post processing site. At the beginning of the PP Phase, the originating site will decide to broadcast either the update set or the relations that haven't been broadcast. In the latter case, the update set can

be constructed at any site. Of course, consistency checks must be performed before the corresponding relations are updated. The query is completed after this phase.

An example of query processing in DDHMIN is illustrated in Figure 3.4, 3.5, and 3.6. The query is a relational query with the following predicate.

$$(R_1.a = R_3.a) \text{ AND } (R_1.b = R_2.b = R_4.b) \text{ AND } (R_2.c = R_3.c = R_4.c)$$

The distribution of the original-relations and the resulting site-relations are shown in Figure 3.4. The various phases of query processing are shown in Figure 3.5.

3.4.2 Improvements to The Five-Phase Query Processing

Because of the expansion effect of joins, the size of a site-relation at a site may be much larger than the total size of the original-relations there. To reduce the cost of constructing site-relations, site-relations are not really constructed in LP phase. Instead, original-relations are used. Site-attributes are extracted from the original-attribute to broadcast only when they are needed. As a result, minor changes have to be made in several phases. In the LP Phase, only local semi-joins, instead of full joins, are performed in each related site among the original-attributes. Since only equi-joins are considered, the set of distinct values in a site-attribute at a related site is the same as that in any original-attribute of the same domain there. Therefore, broadcasting a site-attribute R_x^j in the GSJ Phase is equivalent to broadcasting any original-attribute of domain j in site x . After an attribute is broadcast, all original-relations are semi-joined with this attribute. In the RT Phase, fragments of the original-relations, instead of the site-relations, are broadcast. A unique copy of every relation must be selected to broadcast. Since the directory is redundant and the

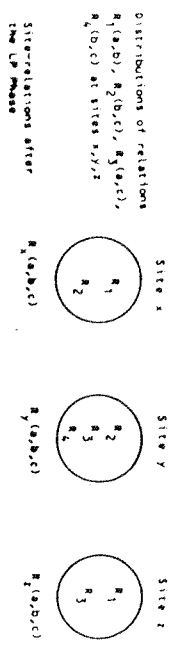


Figure 3.4 Distribution of relations and the result of the LP Phase.

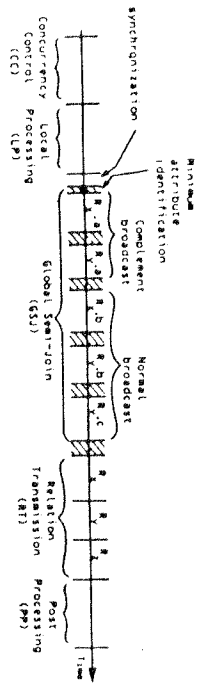


Figure 3.5 An example of the schedule generated by the basic query processing strategy.

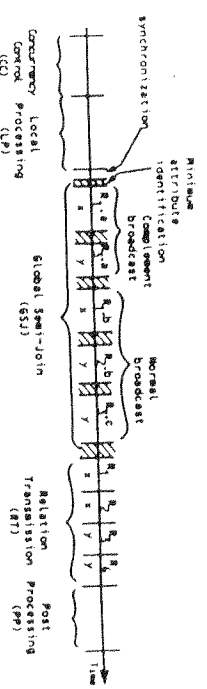


Figure 3.6 An example of the schedule generated by the improved query processing strategy.

size difference among original-relations tends to be small after the GSJ Phase, a tie-breaking rule using the site number can be applied to identify the unique copy to broadcast. An example of the broadcast sequence using the original-attributes for the schedule in Figure 3.5 is illustrated in Figure 3.6.

3.4.3 Protocol for Extremum Identification

In the proposed query processing strategy, it is necessary to synchronize all related sites at the end of each phase and to identify the minimum attribute in each step of GSJ phase. These two problems are equivalent to finding the site with the extremum value and can be solved by the protocol discussed below. Although the proposed protocol does not require additional hardware support, extremum-identification protocols with special hardware interface can also be used [WAH83] [JUA84] [WAH85b]. Without loss of generality, only the problem of identifying the site with the maximum value is studied here. The algorithm is shown in Figure 3.7. It is a distributed algorithm, executed at each site, that determines whether the local parameter, y , is the maximum. Each site contends for the bus and broadcasts its local value. It also listens to the bus for values broadcast by other sites. If the value received is better than y , then it drops out of contention. This process continues until all sites have either broadcast or h have been eliminated. The last broadcast value is then the maximum.

This process is complicated by the non-deterministic time to generate the local parameter at each site. After a site has broadcast its local parameter, it sets a timer to allow sites with larger parameters to broadcast their parameters and sites that have not finished evaluating their parameters to broadcast an 'unfinished' message. If a timeout occurs before such a message is received, it will presume that the current maximum is the real maximum. Since the elapsed


```

process maximum (x, y) {
  /* This process executed at site x with value y identifies whether y is the maximum of all sites.
  current_max represents the broadcast maximum.
  next_x is the site that has been selected to broadcast next */
  current_max := -∞;
  compute next_x; /* estimate the site that will broadcast next */
  fork {generate_y}; /* create process to generate y */
  while (true) do {
    wait (event);
    case (event) of {
      (y is generated): {
        if (x == next_x) then /* broadcast y only if this site is next_x */
          signal network layer to contend for bus and to broadcast y;
        case (y) of {
          (y > current_max): current_max := y;
          (y < current_max): return(failure);
        }
      }
      (receive new_current_max): {
        current_max := max(new_current_max, current_max);
        next_x := x; /* any site could broadcast next */
      }
      (y is not available): (set timer; signal to contend for bus to broadcast 'unfinished');
      (y > current_max): (set timer; signal to contend for bus to broadcast y);
      (y < current_max): {
        if (trying to broadcast y) then withdraw this broadcast;
        return(failure);
      };
      (receive 'unfinished'): {
        unset timer; /* there are unfinished sites */;
        next_x := site sending the 'unfinished' message;
        if (trying to broadcast) then withdraw this broadcast;
      };
      (get bus): {
        if (y is available) then {broadcast(y); set timer}
        else {broadcast('unfinished'); unset timer};
      };
      (timeout): {
        if (trying to broadcast) then {
          /* timeout occur, all unfinished sites stop any further broadcast */
          withdraw this broadcast;
          return(failure);
        }
        else return(success);
      };
    }
  }
}

```

Figure 3.7 Protocol to identify the site with the maximum.

time before receiving such a message depends on the network load, the delivery of these messages should be given higher priority, especially for 'unfinished' messages. To prevent the network from flooded by the unfinished messages, all sites that are still evaluating their parameters will refrain from sending further 'unfinished' messages once the first 'unfinished' message is received. The sequence of events is illustrated in Figure 3.8. The shaded interval indicates a contention period during which sites with larger parameters or 'unfinished' messages contend for the bus. The contention period ends when an 'unfinished' message is broadcast.

The performance of the proposed protocol can be enhanced by making a good initial estimate of the site containing the maximum. In the optimal case, only one broadcast and a timeout period are necessary. If the distribution of the parameters is known, then the maximum can be identified in a constant number of messages (on the average) independent of the number of sites [WAH83] [UA84] [WAH85b]. On the other hand, if the initial estimate is random, the average number of messages is $O(\log_2 n)$, where n is the total number of sites.

The above extremum-identification procedure can be applied to identify the completion of the LP Phase. Parameter y in Figure 3.7 represents the local completion time, and the site with the longest completion time is sought. The time that each site spends between the Initial Broadcast and the end of the LP Phase is determined by the length of the CC and LP Phases. Since the time to resolve conflicts in each site is known globally (all remote actions are broadcast), and the processing times of selections, projections, and local semi-joins are proportional to the cardinality of the relations, a good estimate of the slowest site can be made by all sites. All other sites must wait for the estimated slowest site to broadcast its completion time before determining the next 'unfinished'

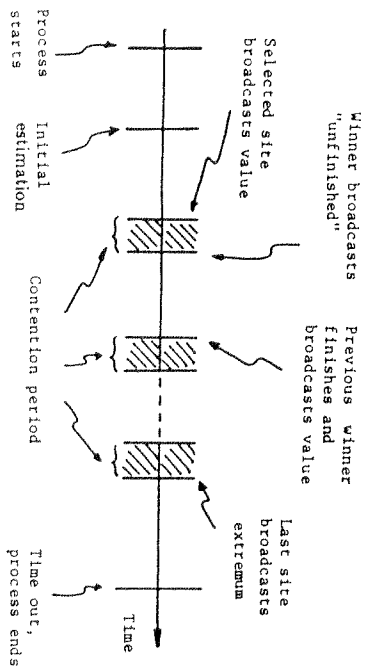


Figure 3.8 The sequence of events in identifying the maximum.

site. To simplify the problem, the time of broadcast can be considered as the completion time of the slowest site, and only a 'finished' message to indicate completion has to be broadcast. After the estimated slowest site has broadcast its 'finished' message, the remaining unfinished sites contend for the bus to elect another site to broadcast the next 'finished' message. If all sites have finished, the last site waits for a timeout period before asserting that the LP Phase has completed.

The procedure in Figure 3.8 can also be applied directly to identify the completion time of each semi-join and the minimum attribute in the GSJ Phase. After one elected site has broadcast its heuristic value, finished sites with larger heuristic values are dropped out of contention, and unfinished sites contend to elect the next site to broadcast. If all sites have finished, sites with smaller heuristic values contend to elect the site with the minimum heuristic value. Since the estimation of the site containing the minimum attribute could be computationally expensive, it may be better to randomly elect it. Another improvement is to allow each site to estimate its heuristic value before it is available, and to drop out of contention if the estimated value is larger than the broadcast value. The next GSJ phase do not have to wait this site. In this way, the slower sites with large heuristic values will not slow down the GSJ phases.

3.4.4 Computing the Heuristic Value

The heuristic function, which estimates the potential of leading to the optimal solution in the GSJ Phase, is another important factor to the system performance. The system performance is affected by the following factors.

- (1) The size of an attribute is important because it affects the broadcast cost directly. Only the *effective size*, which is the product of the domain width

and the larger of the current cardinality of the attribute and the cardinality of its initial complement, should be used.

- (2) Attributes with larger size-reduction effects on other attributes should be broadcast first because they affect the broadcast cost in successive semi-joins.
- (3) The number of unbroadcast attributes on a common domain is another factor of consideration. When this number is large, the broadcast of any attribute in this domain will have a greater size-reduction effect on other relations.

It is very difficult to have a heuristic function that satisfies all the these rules. Depending on the application, the designer has to combine these parameters with appropriate weights to get a good heuristic function.

As defined in section 3.2.2.2, the selectivity of an attribute is usually used to estimate its size-reduction effect, which is defined as the ratio of the current cardinality of the attribute to the initial cardinality of the domain containing the attribute. When an attribute is semi-joined with another attribute, its new size is computed as the product of its current size and the selectivity of the other attribute. This estimation method is not suitable in a broadcast network because the selectivity of a broadcast attribute may be used more than once in computing the size of other attributes. For example, suppose attributes $R_{x,j}$, $R_{y,j}$, $R_{z,j}$ with sizes $s_{x,j}$, $s_{y,j}$, $s_{z,j}$ and selectivities $p_{x,j}$, $p_{y,j}$, $p_{z,j}$ are the attributes in domain j that are broadcast in sequence. After $R_{x,j}$ is broadcast, the sizes of $R_{y,j}$ and $R_{z,j}$ are changed to $s_{y,j} \times p_{x,j}$ and $s_{z,j} \times p_{x,j}$ and their selectivities are changed to $p_{y,j} \times p_{x,j}$ and $p_{z,j} \times p_{x,j}$. When $R_{y,j}$ is broadcast next, the size of $R_{z,j}$ is estimated as $s_{z,j} \times p_{x,j} \times p_{y,j}$ in which $p_{x,j}$ is used twice. This contradicts to the fact that any attributed can only be reduced by another attribute at most once.

To overcome the above problem, the *relative selectivity* is used to measure the size-reduction effect of broadcast attributes. The *relative selectivity* of an attribute in domain j is the ratio of its current cardinality to the cardinality of *current-domain* j , where current-domain j is domain j with the reduction effects taken into account. When an attribute in domain j with relative selectivity p is broadcast, the cardinalities of all attributes in domain j , as well as the cardinality of current-domain j , are reduced by a factor p . As a result, the relative selectivities of all attributes in domain j will not be changed, and the duplicate reduction effects using selectivities will not occur. In the last example, suppose $p_{x,j}$, $p_{y,j}$ and $p_{z,j}$ represent the relative selectivities. Initially, $p_{x,j} = p_{x,j}$, $p_{y,j} = p_{y,j}$ and $p_{z,j} = p_{z,j}$. After $R_{x,j}$ is broadcast, the sizes of $R_{y,j}$ and $R_{z,j}$ are changed to $s_{y,j} \times p_{x,j}$ and $s_{z,j} \times p_{x,j}$ but their relative selectivities remain unchanged. Hence, when $R_{y,j}$ is broadcast, the reduction effects on $R_{z,j}$ is estimated to be $p_{x,j} \times p_{y,j}$.

Four different heuristic functions have been evaluated:

- (1) effective size of an attribute;
- (2) effective-size \times selectivity (selectivity in the conventional sense);
- (3) relative selectivity; and
- (4) effective-size \times relative-selectivity.

These heuristic functions were evaluated by simulations on randomly generated relations. The number of relations and the size of each are limited by the time it needs to find the optimal semi-join schedule in exhaustive enumeration. On the VAX 11/780 computer, it took about three minutes CPU time to simulate a case of three relations with two common attributes each, and 24 hours of CPU time

to simulate a case of three relations with three common attributes each. It was assumed that one copy of each relation existed in the database, and that each existed at a distinct site. The widths of the three attributes were six, eight, and ten respectively. The width of tuples in each relation, including attributes not in the joining domains, was thirty. The domain of each joining attribute was assumed to be integers between one and ten. The distinct values in an attribute were generated by first selecting a threshold, and then including a domain value in the attribute if a random number generated exceeded the threshold. After the attribute values were generated, a cross-product was formed for each relation to obtain the set of possible tuples. Tuples in this set were randomly selected for the relation by a process similar to that of selecting the attribute values. A number of database instances were obtained by using different thresholds.

The simulation results are shown in Table 3.1. The ratios of the total communication costs of various heuristic schedules to that of the optimal schedule are compared. These results demonstrate that dynamic query processing, using relative selectivity as the heuristic function, is very effective. The merits of the heuristic functions must be evaluated again when the relations are generated differently.

3.5 Intelligent Query Scheduling

To take the advantage of available history knowledge in a relatively stable database system, an expert-system approach is proposed here. An expert system is included into the DDBMS to help the transaction handler to semi-join schedules. The interactions between the semi-join scheduler and the *Scheduling Expert System* may be in three different modes: **real-time**, **offline**, and **hybrid**.

Table 3.1
Relative performance of different heuristic schedules as compared to the optimal schedule.

(a) Average behavior of 84 cases for three relations with two common attributes each.

	Effective-Size	Effective-Size	Relative	Effective-Size
	selectivity	selectivity	Selectivity	Relative-Selectivity
Average	1.912	2.037	1.138	2.037
Stand. Dev.	0.338	0.388	0.080	0.388
Maximum	2.537	2.537	1.200	2.537
Minimum	1.099	1.099	1.022	1.099

(b) Average behavior of 17 cases for three relations with three common attributes each.

	Effective-Size	Effective-Size	Relative	Effective-Size
	selectivity	selectivity	Selectivity	Relative-Selectivity
Average	1.742	1.070	1.002	1.058
Stand. dev.	0.459	0.215	0.004	0.206
Maximum	2.200	1.878	1.017	1.878
Minimum	1.000	1.000	1.000	1.000

In the *real-time mode*, the expert system provides the necessary knowledge to the scheduler in the scheduling. Both static and dynamic scheduling are capable of operating in this mode. In the *offline mode*, the expert system provides the scheduler with precompiled schedules stored in its Knowledge Base. *Query-reusable systems* that store precompiled schedules and apply them later fall into this category. Only static scheduling can operate in this mode. Lastly, in the *hybrid mode*, the expert system provides precompiled schedules as well as the necessary expertise to the scheduler. Only dynamic scheduling is suitable in this mode. The expert-system approach allows better scheduling algorithms to be employed that are often too complicated to use in real time. Moreover, the expert system can learn and adjust its scheduling strategies according to changes in the database.

The logical structure of the proposed approach is illustrated in Figure 3.9.

The semijoin scheduler in a Transaction Handler is the part of the code that generates semijoin sequences. When processing a query, the semijoin scheduler consults the Scheduling Expert System to get ready-to-run schedules or scheduling expertise. The Scheduling Expert System has the following modules:

- (1) *Pattern Recognizer*: This module classifies the given queries coupled with the database states into system recognizable classes. Without this module, the input domain of the Scheduling Expert System is too large to manage.
- (2) *Expert Scheduler*: This module is the inference module that serves the Semijoin Scheduler and the Learning Module by returning schedules in its Knowledge Base. In the simplest form, this module can be a knowledge retrieval system.
- (3) *Knowledge Base*: This module stores all the knowledge of the system. The knowledge can be inference rules, expertise of semijoin scheduling,

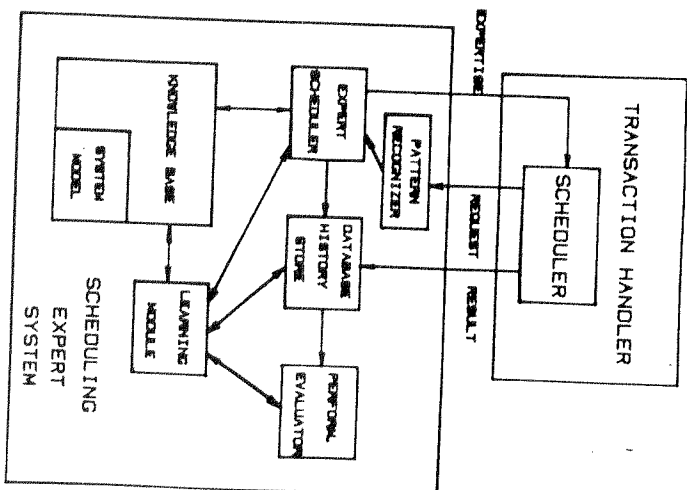


Figure 3.9 Intelligent semijoin scheduling model.

weighting functions, or even the database history. A system model can also be maintained in the Knowledge Base.

- (4) *Database History Store:* The most recent database history including queries with schedules, results, and database states is stored here. This raw knowledge is analyzed by the Learning Module periodically.

- (5) *Performance Evaluator:* When requested by the Learning Module, this module evaluates the performance of the schedules. It is not practical to get an absolute measure on how good a schedule is because the system cannot exhaust different schedules for one query on the same database state. The evaluation may have to be statistical.

- (6) *Learning Module:* This module analyzes the database history to extract the scheduling expertise. It can be as simple as a copy mechanism that merely copies the database history into its Knowledge Base, or it can extract the useful knowledge from the history and compiles them into statistics in the Knowledge Base. It is the heart of the Scheduling Expert System for the following reasons. First, no substantial experience has been collected and built into the database at design time, hence new experience must be learned. Second, the scheduling problem is an NP-hard problem, and no simple rule can guarantee its optimality. More experience would probably help in producing better schedules. Lastly, the performance of the scheduling strategies is application-dependent as well as time-dependent. Thus learning capability can help to adapt to changes in the system.

3.6 Intelligent Scheduling in DDBLMN

The current major difficulty of implementing such an intelligent scheduling system in a DDBMS system is that we do not have enough real-time experience with DQP. We also do not have access to a real DDBMS. Currently, only pure dynamic scheduling is adopted in DDBLMN. This may suffer from high local processing overheads due to possible redundant semijoins in a system with large disk overheads. Further, the communication overheads in synchronization and in identifying the site with the minimum attribute may hamper the effectiveness of dynamic scheduling. To reduce these overheads, we will adopt more intelligent scheduling in DDBLMN in the future.

Three alternatives in implementing intelligent scheduling will be considered. In the first case, dynamic scheduling is used with additional knowledge provided by the Scheduling Expert System to estimate the site with the minimum attribute in semijoins and the slowest site in synchronization. The number of network messages is expected to be reduced, while the local processing overhead will not be reduced significantly. The second alternative, *adaptive hybrid scheduling*, combines static and dynamic scheduling. The scheduler will schedule semijoin broadcasts according to precompiled schedules. When a scheduled broadcast cannot proceed because one or more related relations are locked by other transactions and are inaccessible, the scheduler either chooses the next candidate for broadcast or switches to dynamic scheduling. In this way, both communication and local-processing overheads are expected to be reduced. The cost of implementing the Scheduling Expert System in the above cases would be approximately the same. In the third alternative, scheduling expertise extracted from the database history is used by the scheduler. This is most likely to exist in the form of weighting functions.

The scheduler produces a schedule for each query with the help of the Scheduling Expert System and schedules the semi-join broadcasts as an adaptive hybrid scheduling. The biggest challenge in this approach is the design of the Learning Module.

3.7 Summary

In summary, the most important problem in distributed query processing is to find a strategy to identify the target-data with the minimum processing cost in the processing of multi-relation joins. For a system in which the communication cost dominates the overall processing cost, the most promising approach is to execute a sequence of semi-join to reduce the raw-target before a complete join is executed. A five phase query processing strategy with redundant materialization and dynamic semi-join scheduling is proposed in this chapter. Coupled with concurrency control and redundant materialization, dynamic query processing strategies provide better flexibility than static strategies. Further, the propagation of errors in static semi-join scheduling is avoided. An $O(\log_2 N)$ algorithm based on the characteristic of local multiaccess/broadcast networks is developed to provide an efficient method of collecting statistics in the dynamic strategies.

Finally, an intelligent query processing model is proposed for future research to improve the efficiency of query processing.

CHAPTER IV

CONCURRENCY CONTROL

In this chapter, a concurrency control protocol for local multiaccess/broadcast networks based on the conflict graph and Thomas' [THO79] algorithm is presented. Using the special characteristics of the network, we modify existing protocols of wide-area networks, such that the proposed protocol can suit to this local system. The broadcast capability and short propagation delay time of the network simplify the distribution of complete transaction status to every site. Hence, the protocol behaves similarly as a single site concurrency control protocol.

4.1 Introduction

In Bernstein and Goodman's paper [BER81b], concurrency control is defined as "the activity of coordinating concurrent accesses to a database in a multiuser environment." Concurrency control permits users to share a database while preserving the illusion that each user is executing alone in a dedicated system. The consistency of a DDBMS is endangered by three factors that have been stated in Chapter I.

- (1) Multiple copies of a piece of data may be stored, accessed, and updated by different users at different sites at the same time.
- (2) Communication delays usually prohibit instantaneous distribution of system status information.
- (3) The information exchange may be impossible when some components fail.

The following are two typical anomalies that may be resulted in an ill-controlled electronic funds transfer system.

- (1) **Lost Updates.** Suppose two customers try to deposit money into the same account simultaneously through automatic teller machines (ATM). These two transactions may interfere with each other. The two ATM's handling the two customers could read the account balance at approximately the same time, compute new balances in parallel, and then store the new balances back into the database. One deposit is overwritten by the other and is lost.

- (2) **Inconsistent Retrievals.** Suppose a company has a saving account and a checking account in the system, and each initially contains balances of \$1000. A transaction TR_1 is to transfer \$500 from the saving account to the checking account, and another transaction TR_2 is to query the total balance of both accounts. In the absence of concurrency control, TR_2 may read the balance of the saving account before TR_1 is processed (\$1000) and read the balance of the checking account after TR_1 is processed (\$1500). The total balances that TR_2 gets is \$2500, rather than \$2000.

The objectives of concurrency control are

- (1) correctness,
- (2) maximizing the system throughput,
- (3) minimizing the overhead, and
- (4) simplicity.

Objective (1) is essential in a concurrency control protocol, while a compromise among objectives (2), (3), and (4) may be needed. Simplicity is important in practice since a complicated protocol is hard to prove to be correct, hard to implement, and less robust.

4.1.1 Consistency and Serializability

The correctness of the concurrency control algorithm depends on the requirements of the database. When the processing order of a set of transactions is order independent, i.e. any order is acceptable, the general condition to be enforced is *serializability* [PAP79]. Let E denote an execution of transactions TR_1, TR_2, \dots, TR_n . E is a *serial execution* if no transaction is executed concurrently in E ; that is, each transaction is executed to completion before the next one begins. An execution is *serializable* if it is computationally equivalent to a serial execution, that is, if it produces the same output and has the same effect on the database as some serial execution. The concurrency control of a DBMS should ensure that all executions are serializable.

The *read (or write) set* of a transaction is the set of data items that are read (or written) by the transaction. A *conflict* occurs between two transactions if the write set of one transaction intersects with the read set (r-w conflict) or write set (w-w conflict) of the other. A set of transactions can be executed in parallel and are serializable if the execution order of conflicting transactions is carefully arranged.

4.1.2 Throughput

The system throughput is defined as the average number of transactions a system can process in a time unit. In a real DBMS, the definition is a little vague since the lengths of transactions and the requested data entities are variant. In some cases, the *'degree of concurrency'*, which is the number of transactions that can be processed at the same time, can also be used to measure the performance of the concurrency control protocols. We use *'throughput'* and *'degree of concurrency'* interchangeably in this chapter. More

precise definitions can be found in Papadimitriou's papers [PAP79] [PAP82].

Given a set of transactions, we can always find a serializable execution with the maximum throughput. However, it's not practical to achieve in a real DDBMS for the following reasons.

- (1) The arrival of transactions and the behavior of transactions are difficult or even impossible to predict.
- (2) Even if the transactions are predictable, finding an optimal solution is computationally expensive. To determine whether an execution is serializable or not has been proved to be NP-complete [PAP79].
- (3) An optimal schedule may result in longer delays for some transactions. A real DDBMS may wish transactions to be serviced in a *first-come-first-serve* order.

As a consequence, the most commonly used strategy in concurrency control protocols is a combination of '*first-come-first-serve*' and '*wait-and-error*'. When a transaction enters the system, it is immediately processed until it is terminated or is stopped for some reasons. In the latter case, the stopped transaction either waits until the cause is resolved or restarts.

4.2 Previous Work

Many concurrency control algorithms have been proposed before. They can be classified into three categories: locking, timestamp ordering, and optimistic. Bernstein and Goodman provided a good survey on locking and timestamp ordering based protocols [BER81b]. In a *locking protocol*, data items are locked before they are accessed. Locks must be issued in such a way that the transactions are serializable, and if deadlock occurs, recovery must be possible. Locks are released after the access is completed. The most popular locking

protocol is the *two-phase protocol* in which unlocks for a transaction are issued after all locks have been made. The execution is divided into two phases. During the first phase, called the *growing phase*, the transaction can only request locks, and during the second phase, called the *shrinking phase*, the transaction can only release locks.

In a *timestamp-ordering protocol*, a transaction is assigned a unique timestamp equal to the value of the local clock when it is initiated. In addition, each data item is also timestamped with its most recent read and write. Before a transaction is processed, its timestamp is compared with the timestamps of the accessed data. A transaction is rejected if a conflict is detected. When a transaction is rejected, it should be restarted with a new timestamp, and all intermediate results have to be discarded. An important timestamp ordering protocol is Thomas's *majority consensus timestamp ordering protocol* [THO79].

In an *optimistic protocol* [KUN81] [BHA82b], each transaction works on a private copy of the database and no control is imposed on the execution of transactions. The consistency is checked on the termination of a transaction. The transaction is terminated if it has operated on a consistent state; otherwise, the transaction is restarted.

4.2.1 Comparisons Among the Three Types Of Protocols

In both timestamp-ordering and optimistic protocols, there is no overhead on lock management. A very high degree of concurrency can be achieved if no conflict happens. Obviously, the success of both protocols is highly dependent on the probability of conflicts among transactions, and the throughput tends to decrease as the system load increases. The major problem of using these two types of protocols on a DDBMS is the penalty of restarts, as the query processing

cost is high is restarting a transaction. In addition to the restart overhead, the additional overhead is needed in both protocols for the authorization or the validation of updates.

The expensive restarts can be avoided in lock based protocols. However, they incur the following overheads when the network does not have the broadcast capability.

- (1) *Messages for requesting/releasing locks.* For each data access, three sets of messages are needed: a set of messages for requesting locks on all accessed data, a set of acknowledgements in granting the locks, and a set of messages for releasing the locks. Additional time is also spent in waiting for locks.
- (2) *Extra overhead for deadlock prevention, detection, and recovery.* Since deadlocks may occur in lock based protocols, additional overheads are incurred on deadlock detection and recovery. The degree of concurrency may be degraded when deadlock prevention is enforced.
- (3) *Extra waiting time for overlocking.* A data item may be locked unnecessarily to minimize the possibility of deadlocks and to enforce serializability. Overlocking exists in protocols that request locks initially and in two-phase protocols.

In summary, information exchange is necessary in all three types of protocols. In a lock based protocol, it is done at the operation (read, write) level to ensure serializability. In timestamp ordering protocols, it is done at the update level, while in optimistic protocols, it is done at the transaction level. When the number of writes is much less than the number of reads, the latter two types of protocols are more efficient.

4.3 Concurrency Control in DDJBLMIN

Based on the characteristic of local multiaccess/broadcast networks, a lock based protocol is proposed in this section for its high efficiency and throughput. Recall the characteristics of local multiaccess/broadcast networks stated in Section 1.4.4.

- (1) The system status is almost completely available to all stations by monitoring the activities on the bus.
- (2) Every station gets the information on the bus almost instantaneously, since it is a multi-point configuration.
- (3) All messages arrive at each site in the same order as they are sent, when messages are not lost.

Due to properties (1) and (2), status exchange on bus networks is extremely efficient and fast. It is a natural to use lock based protocols on these networks. Due to the property (3), the transactions are not distinguished from global transactions. Hence, every transaction refers to the same database state.

4.3.1 General Description

It is assumed that the transaction processing model discussed in Section 1.8 is used, that a relation is the basic unit of a data item, and that an access may be either a read or a write but not both.

Broadcast Transaction Table and Precedence Graph

To record the global database state, each site maintains a *Broadcast Transaction Table (BTT)* that records the status of all transactions and the status of relations (locked or unlocked). All Broadcast Transaction Tables are always identical. Based on the information stored in BTT, a *Precedence Graph (PG)* recording the precedence relationship among transactions is also maintained in each site. A transaction is registered into the BTT and PG since its initial broadcast and is retained until its termination. The precedence relationship between two transactions can be added due to a conflict, deleted due to a termination, but cannot be reversed in any case. The serializability of transactions is ensured by this property.

Interaction between transactions and Concurrency-Control subsystem

When a transaction wants to place a lock on a relation, it makes the lock request to the Concurrency-Control subsystem in the local site. The Concurrency-Control subsystem checks the local BTT and PR to determine the serializability condition. The transaction has to wait until serializability is satisfied. With a common bus and assuming that locks received are processed in real time before a lock request is broadcast, only lock requests that will not be rejected by other sites are broadcast. Once a lock is broadcast, the lock is automatically granted to the transaction and is recorded in the the BTT and PG in each site. The newly granted lock will not conflict with other locks that had been granted since the consistency had been checked by the local Concurrency-Control subsystem. If this causes any conflict to the lock requests under processing, then all those conflicting requests must be rejected. Since at most one broadcast can be made at any time on the bus, it is impossible to grant two

conflicting locks simultaneously. In this sense, the common bus acts as a gateway for concurrency control as in a centralized database.

Information exchange

Essentially, information exchange is not needed for concurrency control in DDHMIN. The read set and write set of a transaction are known to all other transactions through the initial broadcast. Moreover, lock and unlock messages can be embedded into the corresponding network activities. Recall the procedure of query processing in Chapter III, the lock and unlock messages can be embedded in each semi-join broadcast and relation broadcast. Although there is no need to explicitly specify the granted locks or released locks in these broadcasts (they can be figured out from the information in BTT), we recommend that the system designer piggybacks this information on the broadcasts to avoid significant overhead to trace the transactions.

The efficiency of information exchange not only saves the communication overhead, but also saves the time waiting for locks and the time of overlocking. The time of overlocking is saved because a lock is released as soon as it is not used by the transaction.

Given a processing order for a set of transactions, the *latest locking time* of a lock on a relation by a transaction is the time that the transaction begins to use (read or write) it. The *earliest unlocking time* is the time that the transaction no longer needs it. The minimum locking time is the difference between these two time instances. Within this time period, a conflict accessing to the relation will introduce inconsistency into the database. Since the information exchange on the network is almost instantaneous, the minimum locking time can be achieved if each transaction issues the lock request at the

right time and responds to lock releasing quickly. This will be discussed later in this section. It must be pointed out that minimum locking does not imply the maximum degree of concurrency which can only be achieved when all the permutations on the transactions are compared.

Precedence order

The precedence order of transactions is governed by the times that a conflict is found. The algorithm detects access conflicts between two transactions from their read and write sets when they are broadcast; however, the order of processing is not imposed until the conflicting access of one of these transactions is made. After the order of processing is defined between TR_1 and TR_2 (say TR_1 precedes TR_2), then all conflicting accesses made by TR_1 must precede accesses of TR_2 to enforce serializability.

This protocol is deadlock free because preceding transactions will never wait for locks from succeeding transactions.

4.3.2 The Locking/Unlocking Time

The correctness and the performance of the protocol are dependent on the time to lock and unlock relations. These are discussed in this section. Depending on the query processing strategy discussed in Chapter III, the time instances of locking and unlocking for read and write operations are different.

Locking time for read-locks

The time instances for a transaction to read-lock a relation is the time that the relation is first read by the transaction. Consistency checks are performed before the relation is locked. For relations that are semi-joined in the LP Phase, they should be locked before the LP Phase begins. For other relations that belong to the read set of the transaction but are not used in the LP Phase, they should be locked once one of their attributes is either broadcast or semi-joined with a broadcast attribute.

Unlocking time for read-locks

Depending on the storage capacity, there are three different time instances to unlock read-locks. In the first case, when the storage capacity is large enough and the access to a relation is read only, a copy of the relation (with only the attributes to be accessed) is made in the working storage when it is read-locked, and the relation is unlocked immediately. An update on the relation by the succeeding transactions will not affect the values obtained by the current transaction.

In the second case, when duplicate copies cannot be made due to storage limitation, then the relation must be read-locked until it is no longer read by the transaction. The locking time in this case may be too long due to the long query processing delay. Thus, the overall system performance may be degraded drastically.

For a system with some extra storage capacity, the following approach can provide much better performance at the cost of some extra storage. When duplicated copies cannot be made, a projection on the attributes that will be used in the global semi-join phase is stored in the working storage. Thus, these

read-only relations can be unlocked at the same time as that in the first case. These projected attributes provide less reductions in the semi-join operations than those provided by the entire relations in the first approach, as indirect semi-joins cannot be used. As a consequence, the query processing overhead may be slightly increased; however the reduction effects of indirect semi-joins are generally small.

Locking and unlocking times of write-locks

A relation that is updated must be write-locked at the beginning of the PP Phase by the first query in the transaction that updates the relation, and released at the end of the PP Phase by the last query in the transaction that updates the relation. For relations that are read and written in the same transaction, the read and write locks must be released together.

4.3.3 Minimum Locking Concurrency Control Protocol

The detailed description of the protocol is given in this section. The distributed concurrency control algorithm is shown in Figure 4.1, 4.2, and 4.3. The precedence relationships among transactions are stored in a precedence graph and are maintained by procedure *consistency* (Figure 4.3). Procedure *precedencecheck* (Figure 4.2) checks whether two transactions conflict with each other.

Process *cc_minimumlock* (Figure 4.1) is triggered by the arrivals of lock requests, the successful contention in getting the bus, the termination of transactions, and the release of locks. Transactions originating from a site communicate with the local process *cc_minimumlock*. Each transaction is a process, which can be in one of the following states: THINKING, SUSPENDED,

```

Process cc_minimumlock (BTT, PG)
/* This process serves the asynchronous arrivals of lock requests, unlock requests, and termination
of transactions. The relations are locked for the minimum amount of time with respect to precedent
transactions. It is assumed that a transaction can lock one relation in each request, and that only
one request is processed at a time (due to the single bus). A transaction that has been initiated
can be in one of the following states: PENDING: all lock requests received must be checked first
to determine if they conflict with the current request; ACTIVE: the transaction is waiting for the
availability of the bus; SUSPENDED: the transaction is not processed because there is another
transaction that is ACTIVE or PENDING, or it conflicts with another transaction that has been
broadcast from another site; and THINKING: transaction is not making any lock request.
Information about the active transactions and locked relations are available in BTT, the Broadcast
Transaction Table. Information about the precedence relations of transactions are stored in PG.
*/
while (true) do {
  wait (event);
  case (event) of {
    (arrival of lock request from TR): {
      if ((there is a PENDING or ACTIVE request) then
        set state of TR, to SUSPENDED and suspend TR,
      else {
        PG = PG; /* make temporary copy of PG */
        if (consistency(TR, PG) = success) then
          If (new lock requests have been received from other sites) then
            set state of TR, to PENDING
          else {
            set state of TR, to ACTIVE;
            signal network layer to contend for bus }
        else {
          if (consistency(TR, PG) = waive) then { update BTT; return (waive) }
          else set state of TR, to SUSPENDED and suspend TR;
          if (there are SUSPENDED requests) then wake up one SUSPENDED request ;}
          (receive lock request from TR, from bus): { /* new lock is granted at another site */
            update BTT and PG;
            if (there is a PENDING or ACTIVE request TRi) then {
              if (precedencecheck(TR, TRi, PG) = reject or waive) then {
                if (TRi is ACTIVE) then
                  withdraw bus-contention request;
                if (TRi is to be rejected) then
                  set state of TR, to SUSPENDED and suspend TR,
                else return (waive);
              }
              if (there are SUSPENDED requests) then wake up one SUSPENDED request ;}
              else if (TRi is PENDING) then {
                set state of TR, to ACTIVE;
                signal network layer to contend for bus ;}
            }
            (a transaction is terminated at another site): update BTT, PG, and PG' (if necessary);
            (unlock is received from bus): update BTT;
            (bus is obtained): {
              broadcast lock request for ACTIVE transaction;
              PG = PG';
              update BTT for locks obtained by the broadcast transaction;
              set state of ACTIVE transaction to THINKING;
              if (there are SUSPENDED requests) then wake up one SUSPENDED request ;}
            }
          }
        }
      }
    }
  }
}

```

Figure 4.1 Process for concurrency control with minimum locking time.

```

procedure precedencecheck (TR, PG);
/* This procedure checks the precedence graph PG to determine whether the lock request by TR, a
   conflict-free transaction, can generate the new lock TR, PG, which is the lock
   transaction that is in the PENDING or ACTIVE state. */
/* check the consistency for r-w conflict */
else if (edge exists from Nold to Nnew in PG and there is r-w or w-w conflict) then
  return (reject);
else if (edge exists from Nold to Nnew and there is a w-w conflict) then
  return (success);
}

```

Figure 4.2 Procedure for checking precedence relation between two transactions.

```

procedure consistency (TR, PG);
/* TR is the requesting transaction. Transaction TRi is requesting to read or write relation R and
   is represented as a node Ni in PG. This procedure checks (and updates if TRi is a new
   transaction) PG for conflicts. An edge will only be added only when one of the conflicting
   relations has been accessed. All transitive arcs exist in PG since the precedence relation is
   transitive. */
/* if Ni does not exist in PG then add a node Ni to PG;
   for (all nodes Nj in PG ≠ Ni) do {
     if (there is no edge between Ni and Nj) then
       if (TRi has r-w conflict with TRj) then
         add an edge from Ni to Nj; add all necessary transitive arcs;
       else if (R is locked by TRj) then
         add an edge from Ni to Nj; add all necessary transitive arcs; return (reject);
       else if (R is unlocked by TRj) then
         add an edge from Ni to Nj; add all necessary transitive arcs; return (reject);
     else if (there is an edge from Ni to Nj) then
       add all necessary transitive arcs; return (reject);
     else if (there is an edge from Nj to Ni) then
       add all necessary transitive arcs; return (reject);
   }
/* if TRi is a read request on R then return (success); /* all r-w conflicts have been checked */
/* check TRi for w-w conflicts if TRi is making a write request on R */
for (all nodes Nj in PG ≠ Ni) do {
  if (there is no edge between Ni and Nj) then
    if (TRi has w-w conflict on R with TRj) then
      /* several alternatives exist here, the strategy used minimizes the number of updates */
      add an edge from Ni to Nj; add all necessary transitive arcs; return (reject);
    else if (edge exists from TRi to TRj) and (TRi has w-w conflict with TRj) then
      return (reject); /* write on R for TRi will be waived later if write has not started */
    else if (edge exists from TRj to TRi) and (TRi has w-w conflict with TRj) then
      return (waive); /* R has been modified by another transaction of lower precedence */
  return (success); /* all w-w conflicts have been checked */
}

```

Figure 4.3 Procedure for maintaining the precedence graph.

PENDING, or ACTIVE. When a transaction is in the THINKING state, it does not generate any request to the database. After a transaction, say TR_i , makes a request, consistency checks are made on a temporary copy of the precedence graph to determine if any conflict exists between TR_i and other transactions in the precedence graph. If a conflict is found, then TR_i is SUSPENDED and processed later. When no conflict exists, TR_i is put into the PENDING state, so lock messages received from other sites can be processed. TR_i will be SUSPENDED if a conflict is found in this state; otherwise, it will signal the network layer to contend for the bus and will be put into the ACTIVE state. The ACTIVE state is terminated when the bus is granted to TR_i and the request is serviced. At this time, the original precedence graph is updated: If a lock request received conflicts with TR_i , when it is ACTIVE, then TR_i is SUSPENDED and the request to contend for the bus is withdrawn.

Procedure consistency (Figure 4.3) maintains the precedence graph. It defines the processing order of two transactions if they have (r-w or w-w) conflict on accessing a common relation. Once this order is defined, all accesses with a r-w conflict on a common relation by a transaction of lower precedence must be initiated after the conflicting relation has been unlocked by transactions of higher precedence. For transactions with a w-w conflict, the transaction with a higher precedence can be waived because the effect on the database is equivalent to the update by the other transaction alone. To allow transactions to be processed efficiently, the precedence order should be defined as late as possible because the transactions are processed asynchronously, and the transaction which first makes the conflicting access is unknown at the Initial Broadcast Time. If a conflict is detected when a transaction, say TR_i , is requesting an access, and if an edge does not exist between the two nodes representing TR_i ,

and the conflicting transaction, say TR_1 , in the precedence graph, then an edge is added between the nodes. Transitive edges are also added to the precedence graph since the precedence relationship is transitive. The rules for determining whether a transition is accepted, rejected, or waived are specified in procedure *consistency*.

An example of concurrent accesses of four transactions as controlled by the proposed algorithm is shown in Figure 4.4. Boxed requests are not executed because they are either rejected or waived. In Step 4, when TR_2 requests a write-lock on A, there is a r-w conflict between TR_1 and TR_2 , and TR_2 is delayed because an edge exist in the precedence graph from N_1 to N_2 . In Step 6, the write request for TR_2 can proceed because A has been unlocked. In Step 8, N_3 is added to the graph. Since B was accessed by TR_2 , an arc is included from N_2 to N_3 and a transitive arc from N_1 to N_3 is also added. TR_3 is allowed to proceed because B has been released by TR_2 . In Step 14, the request to update C by TR_2 is waived because the update on C by TR_3 has been carried out and would have over-written the update by TR_2 (an edge exists from N_2 to N_3).

The serializability and deadlock-free properties of the proposed algorithm are proved in the following theorems.

Theorem 4.1: The sequence of requests as controlled by procedure *cc_minimumlock* represents a serializable execution of the transactions.

Proof: Serializability for a read-only write-only model is satisfied when the requests are scheduled according to an acyclic precedence graph [PAF79] [ULL82]. To prove that the sequence of requests as scheduled by *cc_minimumlock* is serializable, it is necessary to show that the precedence graph as maintained by *consistency* is acyclic, and that *cc_minimumlock* follows this

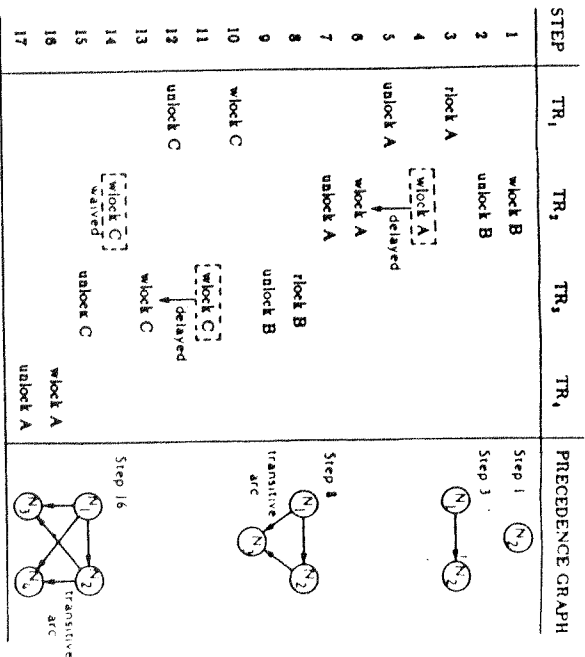


Figure 4.4 An example to illustrate the concurrent execution of four transactions.

precedence graph in scheduling the requests. The precedence graph is maintained by *consistency* is acyclic because transitive arcs are always added whenever a new arc is added to the graph, and no arc will be added between two nodes if an arc already exists between them. To show that the resulting precedence graph is followed by *cc_minimumlock* in scheduling all requests, it can be seen that the precedence relationship of currently active transactions, as defined by the precedence graph is always followed whenever a lock is requested, that an arc is added to the precedence graph whenever a conflicting access is first discovered between two transactions, and that this arc is maintained until one of these transactions is terminated. Hence, the resulting precedence graph obtained at the end of a sequence of requests can be assumed to exist before the requests are executed, and a linear order of execution can be enforced for all transactions, which results in a serializable execution. \square

Theorem 4.2: The proposed protocol *cc_minimumlock* is deadlock-free.

Proof: Given a set of transactions TR_1, TR_2, \dots, TR_n , the wait-for relationship in an execution controlled by *cc_minimumlock* can be represented in the following wait-for graph. For each transaction TR_i , there is a node N_i in the wait-for graph. A directional arc connects N_i to N_j if TR_i waits for TR_j on a conflicting access. Since only transactions of lower precedence must wait for transactions of higher precedence, and the precedence graph is acyclic, the wait-for graph is also acyclic, and the algorithm is deadlock free. \square

The proposed concurrency control algorithm does not impose the precedence relationship between two transactions until the first conflict between them is detected. This allows shorter transactions to proceed first without waiting for the longer ones. However, the maintenance of the precedence graph

at each site may pose a significant overhead. In this case, a predefined precedence relationship may be used when the transactions are initiated. One convenient way is to define a higher precedence for a transaction with a smaller Initial Broadcast Time. The serializability and deadlock-free properties will still be satisfied.

4.4 Local Queries

We have assumed that the status of all transactions is known to all sites in the proposed concurrency control algorithm. However, a local query that only accesses local relations may block other global transactions and this fact is unknown to other sites. To solve this problem, additional status messages can be broadcast during the processing of local queries. This overhead may be large because a database with a good data distribution would have many local queries. Another method is to allow the local queries to be processed when the accessed relations are not locked. When a global transaction TR_1 wishes to access a relation that is locked by a local query TR_2 , it will be blocked in the LP Phase until the relation is unlocked. The site at which TR_2 originates will broadcast an 'unfinished' message when other sites involved in TR_1 are trying to identify the completion of the LP Phase. When TR_2 is finished, this site will broadcast a 'finished' message and start the GSJ Phase for TR_1 .

4.5 Summary

In this chapter, the problem of concurrency control in DDBMS's is examined. Based on the property that information exchange on local multiaccess/broadcast networks is extremely efficient and fast, a lock based protocol using complete system information is developed. The complete system

information is available to each site without additional overhead by monitoring the activities on the networks. A lock based protocol on such networks avoids the large overhead and restarts of a similar protocol in general networks. Although the maximum degree of concurrency is difficult to achieve, the protocol places locks on relations with the minimum locking time. The execution of transactions as controlled by the protocol is proved to be serializable and deadlock free. Through a careful design of locking and unlocking times, the concurrency control protocol is integrated into the distributed query processing strategies.

CHAPTER V

GENERALIZED QUERY PROCESSING

5.1 DQP in Future Communication Networks

Currently, all DQP strategies are developed based on the assumption that the communication cost is in the form of $c_0 + c_1 \cdot x$ where c_0, c_1 are constants and x is the volume of data to be transmitted. Further, the fixed cost c_0 is a cost representing the overhead in setting up the communication channel, and is usually ignored. The objective of DQP becomes the minimization of the volume of transmitted data. This assumption may not be true or even close to the characteristic of future computer networks. In future local networks, data are likely to be transmitted in a store-and-forward fashion rather than through a shared bus. The transmission cost of each "hop" can be modeled as $c_0^{i,j} + c_1^{i,j} \cdot x$, where i and j are the end nodes of the hop, $c_0^{i,j}$ is the fixed cost for a transmission from node i to node j , and $c_1^{i,j}$ is the transmission cost per unit of data from node i to node j . The total communication cost for a transmission along the path $Node_1, Node_2, \dots, Node_n$ can be modeled as

$$\sum_{i=1}^{n-1} (c_0^{i,i+1} + c_1^{i,i+1} \cdot x). \quad (5.1)$$

With high-speed communication technologies, such as optical fiber communication technology, the per unit cost of transmission between each pair

of nodes may be very low as compared to the cost of setting up a communication link, which involves circuit establishment, buffer reservation, queuing, etc. The fixed-link initiation cost is then dominating the communication cost. In this case, Equation 5.1 can be simplified as:

$$\sum_{i=1}^{n-1} (c_0^{i,i+1}). \quad (5.2)$$

Eq. 5.2 shows that the communication cost between any pair of nodes is only a fixed cost for each transmission session. All the DQP strategies being developed for solving multi-site joins, including semi-joins, are no longer adequate. New processing strategies are needed in this case. In this research, we propose a promising strategy as follows:

Algorithm DQPCTSP

- (1) Set current site to the home site.
- (2) Join all queried relations in the current site.
- (3) Send the result to a site storing at least one queried relation that has not been joined. Set this site to be the next current site.
- (4) Repeat steps 2, 3, and 4 until all queried relations are joined together.
- (5) Send the result back to the home site.

An example is shown in Figure 5.1. A query entering site X wants to join R_1, R_2, R_3 , and R_4 together. R_1 in site X is first sent to site Y to join with R_2 and R_3 . The result is then sent to site Z to join with R_4 . Finally, the result is sent back to site X, and the execution terminates. The total communication cost is the summation of the cost incurred in all communication links used in algorithm DQPCTSP. We observe

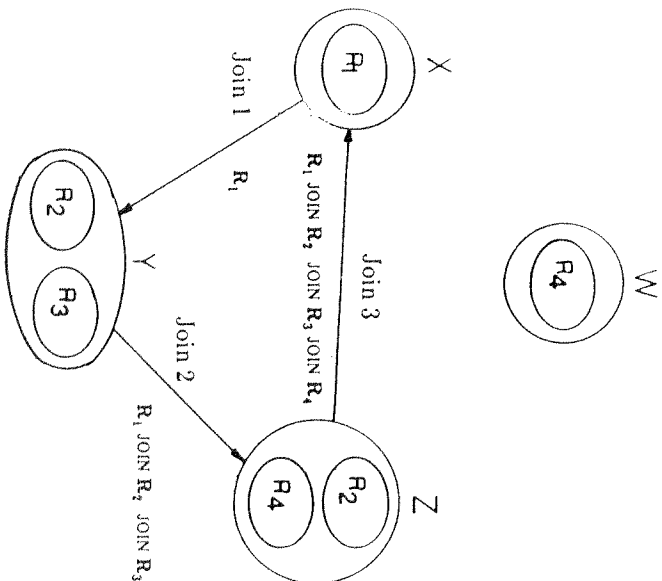


Figure 5.1 An example of Algorithm DQP GTSP.

- (1) that each relation is accessed at least once;
- (2) that each site is activated at most once except the home site; and
- (3) that all the links used by the algorithm connect together as a closed loop.

To minimize the communication cost of Algorithm DQP GTSP, a sequence of connected links that access each relation at least once with the minimum cost should be found. If we cluster all sites storing a given file together, the problem can be viewed as finding a closed path connecting all clusters with the minimum cost. The example in Figure 5.2 shows the result of the clustering on the example in Figure 5.1. Cluster R_1 consists of sites X and Y; Cluster R_2 consists of sites Y and Z; Cluster R_3 consists of site Y only; and Cluster R_4 consists of sites W and Z. The problem can be modeled as the *generalized traveling salesman problems* (GTSP), which is described and solved in the rest of this chapter.

5.2 Introduction to GTSP

The *standard traveling salesman problem* (TSP) can be stated as follows: given a set of cities and the distances (costs) between them, find a tour (a closed path passing through each city exactly once) with the minimum cost [DAN54]. In the *generalized traveling salesman problem* (GTSP), the cities are partitioned into possibly intersecting clusters, and the objective is to find a tour that visits at least one node in each cluster with the minimum cost [SRI70] [KUM70].

Being applicable to solve many real world problems, the TSP have been studied for many years and many efficient heuristic algorithms have been developed. However, the structures of large real world problems tend to be hierarchical. In this case, GTSP will be a better model to represent them.

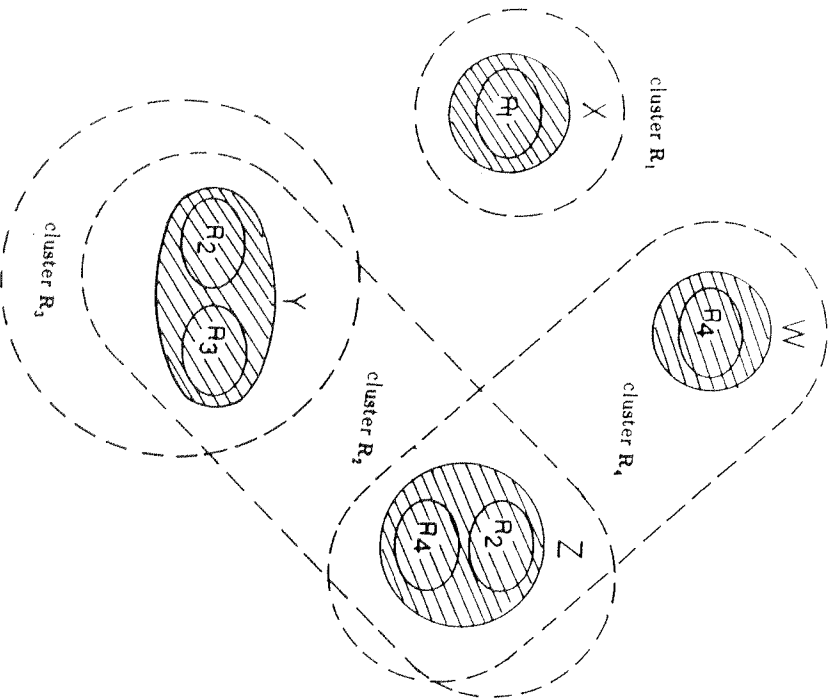


Figure 5.2 The clustering presentation of Figure 5.1.

For example, a traveling salesman wants to visit all dealers in each state of the country. Only one out of a number of possible cities in each state is chosen as the place to meet all local dealers in the state. To minimize the travel expenses, a route visiting a city of every state is to be chosen. This problem can be represented by a GTSP. As a second example, the post office wants to choose the locations of mailboxes for all residential areas (one mailbox for each residential area) as well as a route to collect mail from these mailboxes. A finite number of possible locations in each residential area are under consideration. This problem can also be represented as a GTSP. The third example is the construction of a ring type optical fiber communication network. One of several possible sites is to be selected as the concentrators for each service region. The network is to connect all the concentrators together as a ring. The cost of the construction is proportional to the total length of the network and is to be minimized. The fourth example is a job sequencing problem in which some subtasks are to be cascaded in various ways to carry out a given task with the minimum cost. Each subtask can be carried out by various processes, and the cost of setting up a process depends on the preceding process. By viewing a subtask as a cluster, the process of carrying out a subtask as a city, and the cost of changing from one process to another as the distance between the two cities, the problem can be represented as a GTSP problem. The last example is a distributed query processing problem solved by the algorithm DQP GTSP in Section 5.1. As shown in Figure 5.2, each site storing a copy of queried file is a city; all sites storing a copy of a particular file are clustered into a cluster corresponding to that file; and the traveling cost between a pair of cities is the communication cost of the link connecting the pair of corresponding sites. The algorithm DQP GTSP becomes finding a tour that visits each cluster at least once with the minimum cost.

The problem instances of both GTSP and TSP can be represented as graphs. For any city in a problem instance, there is a node in the corresponding graph representation. The cost of traveling from one city to another is represented by a weighted arc connecting the nodes that represent the cities. Various approaches solving the GTSP are discussed in Section 5.3. Formal graph definition of the GTSP and some properties regarding its optimality are discussed in Section 5.4. A set of transformations solving the GTSP by the TSP-transformation approach are presented in Section 5.5. In Section 5.6, the properties of the transformed graph and the application of existing TSP algorithms to solve the GTSP are discussed. A set of GTSP instances are tested using the proposed algorithm. Finally, a summary is given in Section 5.7.

5.3 Solution Methods of GTSP

The TSP is a well known NP hard problem. It is easy to prove that the GTSP is also NP hard since the TSP is a special case of the GTSP. Three different possible approaches to solve the GTSP are discussed here. In the first approach, the *direct-solution* approach, solution algorithms are directly applied on the given problem instances. There are some standard techniques applicable for solving the problem optimally, such as dynamic programming [SR170] [KUM70], and branch-and-bound algorithm [LIT63] [GUP78]. Heuristic algorithms that give suboptimal solutions are needed for large problem instances. However, the development of efficient and provable heuristics may involve considerable research efforts, since the problem is much more complex than the TSP, which has drawn extensive research for the past 30 years. More extensive research is needed before satisfactory solution algorithms can be

found. The following two approaches utilize existing results in the TSP to solve the GTSP. They may not be better than the direct-solution approach, but a number of available solution algorithms for the TSP can be immediately applied.

In the second approach, the *cartesian-decomposition* approach, a TSP solution is found for each member of the Cartesian product of the clusters. Each TSP solution is a feasible solution of the GTSP. Further, if the graph satisfies triangle inequality, then the TSP solution with the minimum cost provides an optimal solution to the GTSP. Assuming that there are m clusters with k_i nodes in cluster i , the total number of TSP's to be solved is $\prod_{i=1}^m k_i$. In the third

approach, the *TSP-transformation* approach, the original problem is transformed into another problem such that the TSP solution to the transformed problem provides a solution to the GTSP of the original problem. In this chapter, we show a transformation procedure in which every node in a cluster of the original graph is expanded into three nodes. A single TSP with $\sum_{i=1}^m k_i + m$ nodes has to be solved.

In general, the Cartesian-decomposition approach has less complexity than the TSP-transformation approach when exponential-time algorithms are used to solve the standard TSP. On the other hand, the TSP-transformation approach is more efficient than the Cartesian-decomposition approach if polynomial-time approximation algorithms are used.

5.4 Graph Representation and Some Optimality Properties

Let $G = (N, A)$ be a directed graph. The set of nodes N are indexed by $0, 1, 2, \dots, n$. For each arc $(i, j) \in A$, there is a non-negative integer $d(i, j)$ representing a cost measure of the arc. The costs satisfy the *triangle inequality*; that is, $d(i, j) + d(j, k) \geq d(i, k)$ for all $i, j, k \in \{0, \dots, n\}$. Nodes $1, 2, \dots, n$ are partitioned into m possibly intersecting subsets (*clusters*) C_1, C_2, \dots, C_m , with $k_i = |C_i| > 0$ nodes in subset C_i , and $k_1 + k_2 + \dots + k_m \geq n$. For convenience, denote $\{0\}$ as cluster C_0 , $C = \{1, 2, \dots, m\}$, and $C^+ = \{0, 1, \dots, m\}$. An arc (i, j) is an *intercluster arc* if i and j are not in the same cluster; otherwise, it is an *intracluster arc*. A sequence of nodes $i_0, i_1, \dots, i_p, i_{p+1}$ is called a *generalized tour (g-tour)* on G , when

- (i) $i_0 = i_{p+1} = 0$,
- (ii) $(i_k, i_{k+1}) \in A$, for $k = 0, 1, \dots, p$, and
- (iii) $\{i_1, i_2, \dots, i_p\} \cap C_k \neq \emptyset$, $k = 1, 2, \dots, m$.

The cost Z of a tour is defined as

$$Z = \sum_{k=0}^p d(i_k, i_{k+1})$$

The GTSP is to find a minimum-cost g-tour (optimal g-tour) on the graph G over all such g-tours on G . If $k_i = 1$, for $i = 1, \dots, n$, and $C_i \cap C_j = \emptyset$, for $i, j \in \{1, \dots, n\}$, $i \neq j$, then the GTSP degenerates into a TSP.

Call a cluster that is (not) intersected with any other cluster in a graph as a *intersecting (non-intersecting) cluster*, and define the *covering cluster set* of a node as the set of clusters containing the node. We denote the covering cluster set of v_i as $S_i = \{l : v_l \in C_l\}$. Two nodes are *mutually exclusive* to each other (referring to as *mutually exclusive nodes*) if the covering cluster set of one node is

contained in that of the other. A g-tour *dominates* another g-tour if it visits fewer nodes with a cost no more than that of the other. A g-tour is *promising* if it is not dominated by any other g-tour on the graph, and is *non-promising* otherwise. By first excluding as many non-promising g-tours as possible out of the search space in a solution algorithm, the complexity will be greatly reduced. It is not practical trying to identify all such g-tours, because the problem is equivalent to that of finding all optimal g-tours connecting the minimum number of cities. From the definition of promising g-tours, we know there exists at least one promising optimal g-tour in each graph. Therefore, excluding all non-promising g-tours from the solution space will not eliminate all optimal g-tours, although an optimal g-tour may be non-promising. The following definitions and Theorem 5.1 are helpful in identifying some non-promising g-tours in terms of topological properties. A *canonic g-tour* is a g-tour visiting no mutually exclusive node. For example in Figure 5.3, C_1 is a non-intersecting cluster, while C_2 and C_3 intersect to each other. The covering cluster sets of v_1, v_2 , and v_3 are $\{1\}$, $\{2\}$, and $\{3\}$, respectively, while the covering cluster set of v_4 is $\{2, 3\}$. The node v_4 is exclusive to both v_2 and v_3 . G-tour $\{v_1, v_2, v_4, v_3, v_1\}$ is dominated by g-tour $\{v_1, v_2, v_3, v_1\}$ and hence is non-promising. Further, it is not a canonic g-tour, because it visits mutually exclusive nodes, $\{v_2, v_4\}$ and $\{v_3, v_4\}$. The g-tour $\{v_1, v_2, v_3, v_1\}$ is a canonic g-tour.

Since all nodes in a non-intersecting cluster are mutually exclusive to each other, a canonic g-tour visits such a cluster only once. A canonic g-tour on a graph without intersecting cluster is always promising, but it is not true if the graph contains intersecting clusters. Theorem 5.1 shows that a non-canonic g-tour is non-promising. For this reason, we are studying canonic g-tours only in this chapter.

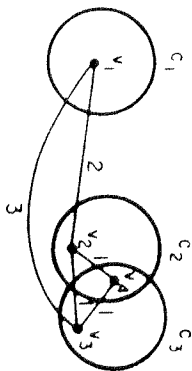


Figure 5.3 Some g-tours on a graph.

Theorem 5.1. For a graph satisfying the triangle inequality, any g-tour that visits at least two mutually exclusive nodes is dominated by another g-tour on the graph.

Proof: Let T be a g-tour that visits two mutually exclusive nodes v_i and v_j , with $S_g \subseteq S_i$. We can construct another g-tour T' with no additional cost which visits the same set of nodes visited by T except v_j . With no loss of generality, let $T = \{t_1, v_1, v_2, v_2, t_2, v_3, v_4, v_4\}$, as shown in Figure 5.4, where v_1 and v_2 (resp. v_3 and v_4) are the nodes connected to node v_2 (resp. v_3), and t_1 (resp. t_2) is the path connecting v_1 and v_4 (resp. v_2 and v_3). A closed path T' ($t_1, v_1, v_2, t_2, v_3, v_4$) as shown in Figure 5.5 can be constructed. T' is a g-tour since all clusters visited by T are visited by T' . The cost of T' is no more than that of T , since $d(v_3, v_4) \leq d(v_3, v_2) + d(v_2, v_4)$. Therefore, T is dominated by T' . \square

The following corollary states that a g-tour visiting more than one node of a non-intersecting cluster is non-promising.

Corollary 5.1. For a graph satisfying the triangle inequality, any g-tour that visits more than one node of a non-intersecting cluster is dominated by another g-tour that visits only one node of the cluster.

Proof: Since the covering cluster sets of all nodes in a non-intersecting cluster are all the same, by Theorem 5.1, such a g-tour is non-promising. \square

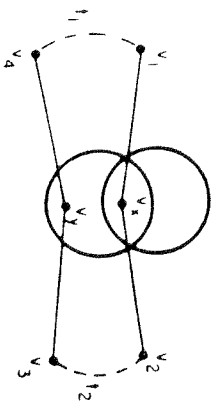


Figure 5.4 A g-tour T containing mutually exclusive nodes (v_1 and v_2)

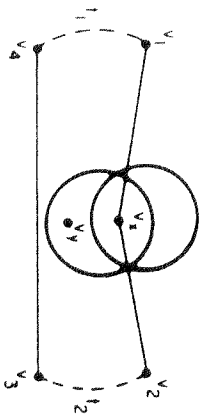


Figure 5.5 A g-tour T' that dominates g-tour T in Fig. 5.4.

5.5 TSP-Transformation Approach

In this section, the GTSP is solved by the TSP-transformation approach. A transformation for graphs with non-intersecting clusters is presented in Section 5.5.1. Another set of transformations are developed in Section 5.5.2 to solve the GTSP on graphs with intersecting clusters. In Section 5.6, heuristic algorithms for TSP, which are applicable to the transformed problem, are discussed and are evaluated on some test cases.

5.5.1 G-tours on Graphs with Non-intersecting Clusters

We develop a G - S transformation that transforms a graph G with non-intersecting clusters into another graph $G' = (N', A')$, such that any tour on G' is a canonic g-tour on G with the same cost and vice-versa. Denote the set of nodes in C_i by $v_{i,0}, v_{i,1}, \dots, v_{i,k_i-1}$ for $i = 1, \dots, m$. The set of nodes N' in G' is obtained as follows:

- (i) For each node $v_{i,j}$ in G , create three nodes $i_{i,j}$ (i-node), $b_{i,j}$ (b-node) and $o_{i,j}$ (o-node) in G' .
- (ii) For node θ in G , create two nodes $i_{\theta,0}$ and $o_{\theta,0}$ in G' .
- (iii) For each cluster C_i in G , create a node q_i (g-node) in G' .

Let C'_i denote the set of nodes $\{i_{i,0}, b_{i,0}, \dots, i_{i,k_i-1}, b_{i,k_i-1}, o_{i,k_i-1}, q_i\}$ in G' created for the set of nodes in C_i . $C'' = \{C'_i, i \in C\}$ are the clusters in G' . We also define the set $\{i_{\theta,0}, o_{\theta,0}\}$ as C'_θ , and $C''' = C'' \cup C'_\theta$. The set A' contains the following arcs:

- (iv) For the i-nodes, b-nodes and o-nodes in cluster C'_i , they are connected into a cycle with arcs: $\{(o_{i,j}, b_{i,j}), (b_{i,j}, i_{i,j}), (i_{i,j}, o_{i,(j+1) \bmod k_i}) : j=0, 1, \dots, k_i-1\}$. These arcs have zero cost.

- (v) The g -node in cluster G_i' is connected by arcs: $\{(i_j, g_i), (g_i, b_{i,j}) : j = 0, 1, \dots, k_i - 1\}$. Additional arcs are created between the b -nodes and o -nodes: $\{(b_{i,j}, o_{i,j}) : j = 0, 1, \dots, k_i - 1\}$. These arcs have zero cost.

- (vi) An arc $(i_{0,0}, o_{0,0})$ with zero cost is created for cluster 0.

- (vii) For an intercluster arc $(i_{i_1, j_1}, i_{i_2, j_2}) \in A$, $0 \leq i_1, i_2 \leq m$, $i_1 \neq i_2$, $0 \leq j_1 \leq k_{i_1} - 1$, $0 \leq j_2 \leq k_{i_2} - 1$, there is a corresponding intercluster arc $(o_{i_1, j_1}, i_{i_2, j_2})$ in A' with the same cost. That is, a cluster can only be entered through i -nodes and departed through o -nodes.

An example graph G is shown in Figure 5.6. There are 24 g -tours. The optimal g tour is $(0, v_{1,1}, v_{2,0}, v_{3,0}, 0)$ with cost 12. The arcs added to G' in Steps (iv), (v) and (vi) are shown in Figure 5.7. Actually, single-node clusters that contain only one node, need not to be expanded.

The costs of intercluster arcs created by step (vii) are the same as the corresponding arcs in G , while the other arcs are intracuster arcs with zero costs. Therefore, the cost of tours on G' consists of the costs of intercluster arcs only. Any tour on G' can only visit clusters in the following restricted way: when visiting a cluster, a tour first visits an arbitrary i -node $i_{i,j}$, then visits all other nodes in the cluster, and leaves the cluster through the o -node $o_{i,j}$. If intracuster arcs are taken away from a tour, then the remaining arcs represent a canonic g -tour on G . This is shown in the following theorems.

Theorem 5.2: Every tour on G' visits the nodes in G_i' , $i \in \{1, 2, \dots, m\}$, using a subtour of the form: $(i_j, o_{i,(j+1) \bmod k_i}, b_{i,(j+1) \bmod k_i}, i_{i,j})$.

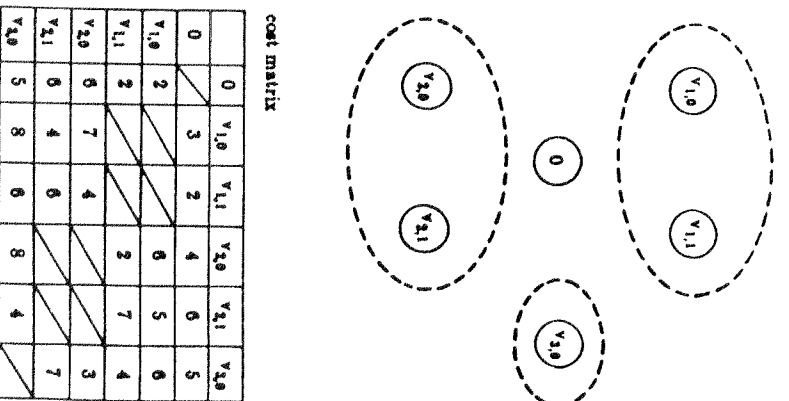


Figure 5.6 Clustered graph G with non-intersecting clusters.

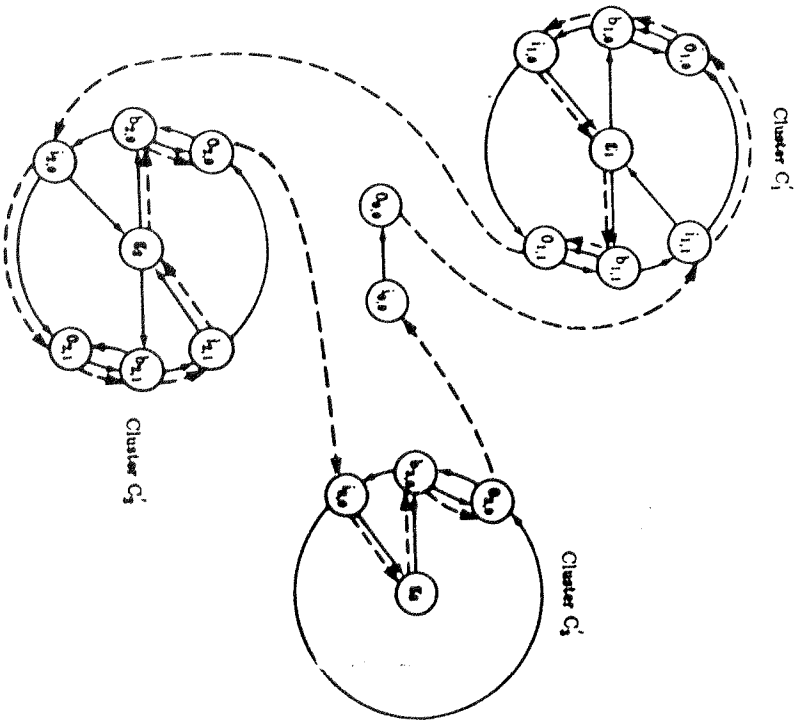


Figure 5.7 Arcs added to G' in Steps (iv), (v) and (vi) of the construction process.

$$i_{1,j} \equiv 1 \pmod{k}, \dots, o_{1,j} \equiv k-1 \pmod{k}, \quad b_{1,j} \equiv 1 \pmod{k}, \quad i_{1,j} \equiv k \pmod{k}, \quad g_{1,j} \equiv b_{1,j} \equiv o_{1,j} \pmod{k} \text{ for } j \in \{0, \dots, k-1\}.$$

Proof: If $i_{1,j}$ is the first node in cluster C_j' that is visited in the tour, then the theorem asserts that all the other nodes in cluster C_j' would be visited before the tour exits through $o_{1,j}$. The existence of such a subtour is obvious from the construction of G' . To show that the subtour is unique, suppose the theorem is false, and the tour exits through o_{1,j_2} , $j_2 \neq j_1$, then a contradiction can be derived. Referring to Figure 5.8, this implies that b_{1,j_2} cannot be entered through o_{1,j_2} and must be entered through g_1 . After traversing node o_{1,j_1} , the tour can either exit the cluster or visit node b_{1,j_1} . In the first case, b_{1,j_1} cannot be visited from o_{1,j_1} and must be visited through g_1 . Node g_1 is, thus, visited twice. In the second case, the tour cannot exit without visiting nodes i_{1,j_1} or o_{1,j_1} again. These imply that the tour must exit through o_{1,j_1} and the theorem is proved. \square

Let Σ_m denotes the set of all the permutations of $\{1, 2, \dots, m\}$. The following three theorems prove the equivalence between the tours on G' and the canonic g-tours on G .

Theorem 5.3: For every canonic g-tour on $G = (N, A)$, there is a tour on $G' = (N', A')$ with $Z' = Z$.

Proof: Let $T = (0, v_{1,j_1}^1, v_{1,j_2}^1, \dots, v_{1,j_m}^1, 0)$ be any given canonic g-tour on G where $(i_1^1, i_2^1, \dots, i_m^1) \in \Sigma_m$ and $j^l \in \{0, 1, \dots, k-1\}$ for $l = 1, \dots, m$. A

As an example, the corresponding optimal tour of G' for the graph in Figure 5.6 is shown as dashed lines in Figure 5.7.

5.5.2 G-tours on Graphs with Intersecting Clusters

For graphs with intersecting clusters, it is not easy to obtain a transformation which provides a one-to-one correspondence between canonic g-tours on a given graph and the tours on the transformed graph. The transformation developed in this section can only guarantee that the optimal g-tours on a graph can be obtained from the optimal tours on its transformed graph and vice versa. A tour obtained from the transformed graph may represent a non-canonic or non-feasible g-tour on G . However, another *canonical transformation* can transform such a g-tour into a canonic g-tour. The transformation for the GTSP with intersecting clusters consists of two stages. First, an *I-N transformation* transforms the original graph into an intermediate graph (*I-N graph*) with non-intersecting clusters. Next, the I-N graph is transformed into a third graph (*G-S graph*) by the G-S transformation developed in Section 5.5. The G-S graph is then solved as a TSP. The I-N transformation that transforms the original graph $G = (N, A)$ into the I-N graph $G'' = (N'', A'')$ with non-intersecting clusters $C_0'', C_1'', \dots, C_m''$ is described as follows.

In a g-tour on G , suppose the tour passes through a node v_i whose covering cluster set is $S_i = \{t : v_i \in C_t\}$, i.e., $v_i \in \bigcap_{t \in S_i} C_t$. In G'' , v_i is expanded into $|S_i|$

nodes, referred to as the *expanded nodes* of v_i , one for each cluster in the covering cluster set. A g-tour that visits v_i in G implies that the corresponding g-tour in G'' would visit all the nodes expanded from v_i .

With $N'' = \emptyset$ and $A'' = \emptyset$ initially, the I-N graph $G'' = (N'', A'')$ is constructed as follows:

- (i) For every node v_i in G , it is expanded into $|S_i|$ nodes in G'' . $N'' = N'' \cup \{v_{i,x} : i \in S_i\}$. Notice that $|S_i| \geq 1$. Denote the expanded nodes of v_i as E_i , i.e., $E_i = \{v_{i,x} : v_i \in C_i\}$.
- (ii) For any two nodes $v_{i,x}, v_{j,y}$, $i, j \in S_i$ in G'' , create two arcs $A'' = A'' \cup \{(v_{i,x}, v_{j,y}), (v_{j,y}, v_{i,x})\}$ with $d(v_{i,x}, v_{j,y}) = d(v_{j,y}, v_{i,x}) = 0$.
- (iii) For every arc $(v_i, v_j) \in A$, let $A'' = A'' \cup \{(v_{i,x}, v_{j,y}) : i \in S_i, j \in S_j\}$ and $d(v_{i,x}, v_{j,y}) = d(v_i, v_j)$.

An example is shown in Figure 5.9. Figure 5.9(a) is part of G . Its corresponding part on G'' is shown in Figure 5.9(b).

Essentially, every node in the intersection of multiple clusters is expanded such that there is one expanded node for each cluster containing the node. The cost on the arcs among the expanded nodes is zero. The arcs directed to and emanating from the original node are duplicated for all the expanded nodes.

Define the *internal tour* corresponding to a node v_i to be the path connecting consecutively all the expanded nodes of v_i in G'' . The cost of an internal tour is zero since the cost of traversing the expanded nodes is zero. The follows theorems show the equivalence between the optimal g-tours on G and G'' .

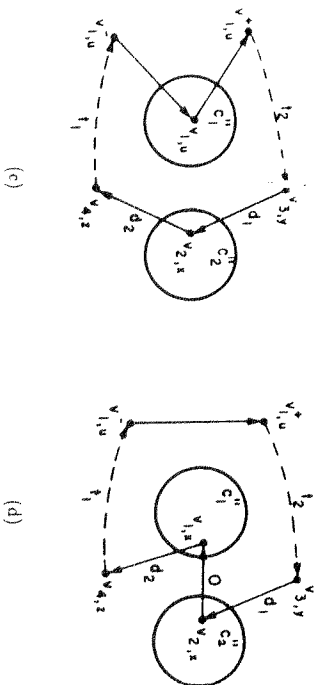
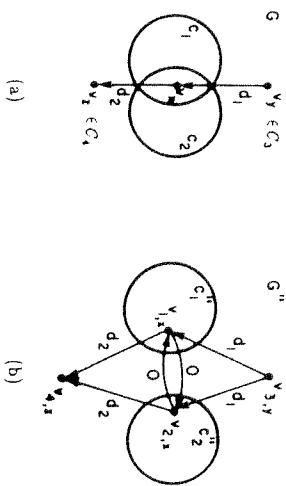


Figure 5.9 An example of I-N transformation;
 (a) a part of graph G ; (b) the I-N graph G'' of G ;
 (c) a g-tour T'' on G'' ; (d) a better g-tour T'' on G'' .

Theorem 5.8. For every g-tour on G , there is a g-tour on the I-N graph G'' with $Z'' = Z$.

Proof: Given a g-tour T on G , we can obtain an equivalent g-tour T'' on G'' as following: For every node v_z visited on G , an internal tour can be created on G'' that visits all nodes $v_{i,x}$ in E_z . For an arc (v_z, v_j) in T , the arc $(v_{i,x}, v_j)$ is contained in T'' , where $v_{i,x}$ is the last node of the internal tour connecting the nodes in E_z , and v_j is the first node of the internal tour connecting the nodes in E_j . It is obvious that T'' is a g-tour of G'' with cost $Z'' = Z$. \square

Though every g-tour on G has an equivalent g-tour on G'' , there is no such correspondence from G'' to G except for the optimal g-tours. This is shown in the next theorem. Denote the preceding (succeeding) node of a node $v_{i,x}$ in a g-tour T'' as $v_{i,x}^-$ ($v_{i,x}^+$).

Theorem 5.7. For every optimal g-tour on the I-N graph G'' , there is an optimal g-tour on the original graph G with $Z = Z''$.

Proof: Let T'' be any promising optimal g-tour in G'' . By Corollary 5.1 in Section 5.4, T'' visits each cluster exactly once. Next, we show that from T'' , we can always construct another optimal g-tour T'' , such that for any E_z in G'' , if T'' visits any node in E_z , then T'' will visit all nodes in E_z consecutively. Without loss of generality, let $T'' = (v_{2,x}^-, v_{2,x}^+, v_{2,x}^-, v_{2,x}^+, v_{1,u}^-, v_{1,u}^+, v_{1,u}^-, v_{1,u}^+, v_{1,x}^-, v_{1,x}^+, v_{1,x}^-, v_{1,x}^+, v_{2,x}^-, v_{2,x}^+, v_{2,x}^-, v_{2,x}^+)$, where l_1 (resp. l_2) is a path connecting $v_{2,x}^+$ and $v_{1,u}^-$ (resp. $v_{1,u}^+$ and $v_{2,x}^-$). T'' visits $v_{2,x}$ in cluster C_2 , but

visits $v_{1,u}$ instead of $v_{1,x}$ in cluster $G_1^{T''}$, where $G_1^{T''}, G_2^{T''} \in S_2$. The node $v_{1,u}$ may or may not be $v_{1,x}$. The g-tour T'' is depicted in Figure 5.9(c). The new g-tour T''' can be constructed by visiting $v_{1,x}$ immediately after visiting $v_{2,x}$.

The arc $(v_{2,x}, v_{2,x}^+)$ is replaced by the arc $(v_{1,x}, v_{2,x}^+)$ and the arcs $(v_{1,u}, v_{1,u}^+)$ and $(v_{1,u}, v_{1,u}^+)$ are replaced by the arc $(v_{1,u}, v_{1,u}^+)$. Since $d(v_{1,u}, v_{1,u}^+) \leq d(v_{1,u}, v_{1,u}^+) + d(v_{1,u}, v_{1,x}^+)$, $d(v_{2,x}, v_{2,x}^+) = 0$, and $d(v_{1,x}, v_{2,x}^+) = d(v_{2,x}, v_{2,x}^+)$, the cost of T''' is no more than that of T'' . Further, since T'' is optimal, the cost of T''' and T'' have to be equal. Therefore, every optimal g-tour T''' on G'' can be reconstructed to another optimal g-tour T'' in G'' which visit all nodes expanded from the same node in G consecutively. From T''' , it can easily be seen that a corresponding g-tour T on G can be constructed (by grouping all the nodes in an internal tour in T''') with the cost $Z = Z''$.

Next, we prove that this g-tour must be optimal on G . Suppose that it were not optimal, then this implies that there is another g-tour on G with cost less than Z . By Theorem 5.6, there exists a g-tour on G'' with cost less than Z . This contradicts the assumption that the g-tour T'' on G'' is an optimal g-tour. \square

From Theorem 5.6 and 5.7, the optimal g-tours on a graph with intersecting clusters can be obtained from the optimal g-tours on its I-N graph. The optimal g-tours on the I-N graph, which contains non-intersecting clusters only, can be obtained from the solution methods developed in Section 5.5.

However, the optimal TSP tours may not be easy to obtain if the problem size is large. We should be able to transform any TSP tour on the G-S graph into a feasible g-tour or a canonic g-tour on the original graph, such that any heuristic TSP solution to the G-S graph provides a heuristic GTSP solution to the original graph. From the theorems in Section 5.5.1, we know every TSP tour

on the G-S graph represents a g-tour on the I-N graph. So we only need to transform each g-tour on the I-N graph into a canonic g-tour on the original graph.

Given a graph G and its I-N graph G'' , an arbitrary g-tour T''' on G'' may visit the expanded nodes of a node v_i of G in the following ways.

- (1) None of the expanded nodes is visited by T''' .
- (2) All expanded nodes are visited by T''' consecutively. That is, they are connected as an internal tour.
- (3) The expanded nodes are visited by T''' in more than one non-consecutive sub-tours. In this case, T''' represents a g-tour T in G visiting v_i more than once.
- (4) Some but not all of the expanded nodes of v_i are not visited by T''' .
- (5) A combination of cases (3) and (4).

Some examples are shown in Figure 5.10 to Figure 5.13. A node v_i in a graph G and its expanded nodes in the I-N graph G'' are shown in Figure 5.10. An example of T''' in G'' , which represents a g-tour visiting v_i more than once (case (3)), is shown in Figure 5.11. Another example of T''' , which visits only part of the expanded nodes of v_i (case (4)), is shown in Figure 5.12. While T''' in Figure 5.13 is an example of case (4), sub-tours t_1, t_2 , and t_3 in each figure are sub-tours. The equivalent sub-tours on G and G'' carry the same symbols.

T''' should be reconstructed except that all expanded nodes are visited by either case (1) or case (2). A canonic transformation algorithm CANTRANS to be described first transforms a g-tour T''' on G'' into another g-tour, which represents a feasible g-tour in G , then transforms any g-tour in G into a canonic g-tour. Given a g-tour T''' in the I-N graph G'' , which is transformed from an

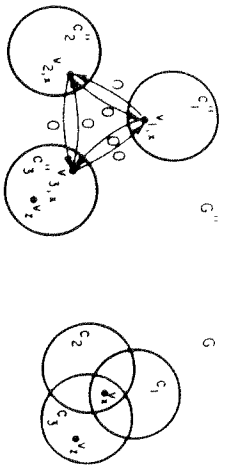


Figure 5.10 A graph G and its LN graph G'' .

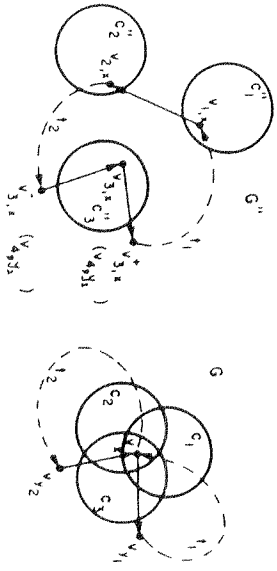


Figure 5.11 A tour on G'' that visits v_2 on G twice.

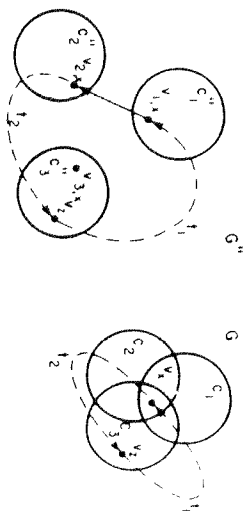


Figure 5.12 A tour on G'' that does not visit all expanded nodes of v_2 .

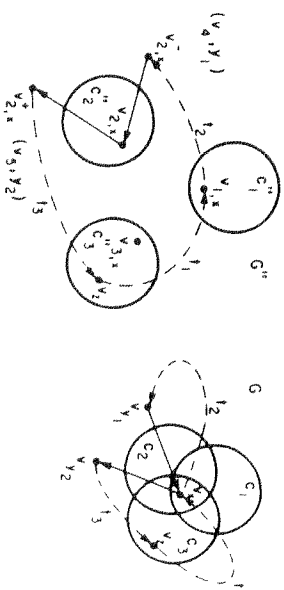


Figure 5.13 A tour on G'' that is a combination of Case (3) and (4).

original graph G , another g-tour in G'' representing a feasible g-tour in G is first reconstructed by applying steps (1) on T'' . Next, the new g-tour is transformed into a canonic g-tour by applying step (2) on every pair of mutual exclusive nodes. The algorithm is shown as follows.

Algorithm CAN-TRANS

(1) For each node v_i with expanded nodes E_i in G'' ,

(1.a) If its expanded nodes in G'' are visited by T'' in either case (1) or case

(2), then do nothing;

otherwise, execute steps (1.b) and (1.c).

(1.b) Assuming that the expanded nodes that are contained in T'' forms several pieces of possibly zero length (a single node) disconnected paths, called *internal paths*,

$(v_{i_1, x}^-, \dots, v_{i_2, x}^+), (v_{i_3, x}^-, \dots, v_{i_4, x}^+), \dots, (v_{i_q, x}^-, \dots, v_{i_{q-1}, x}^+)$, do

following:

Construct the internal tour T_i that connects all nodes in E_i in an arbitrary order except the first node and the last node. There are three different cases in the selection of the first node and the last node of T_i :

Case (a) : If at least one internal path has non-zero length, then select the first node (resp. last node) of such an internal path to be the first node (resp. last node) of T_i . If more than one alternative exists, then select the pair of nodes, say $v_{i_{1, x}^+}$ and $v_{i_{q, x}^-}$, such that the value of $d(v_{i_{1, x}^+}, v_{i_{1, x}^-}) + d(v_{i_{1, x}^+}, v_{i_{1, x}^+})$ is the minimum among all alternatives.

Case (b) : If the lengths of all internal paths are zero, i.e., they are $(v_{i_{1, x}^-}), (v_{i_{1, x}^+}), \dots, (v_{i_{q, x}^-})$, then the selection is the same as Case (a) but

the first node and the last node should belong to two different internal paths. Assuming that the first node is $v_{i_{1, x}^-}$ and the last node is one of $\{v_{i_{q, x}^-}, \dots, v_{i_{q, x}^+}\}$, say $v_{i_{q, x}^+}$, take the arc $(v_{i_{1, x}^-}, v_{i_{q, x}^+})$ to emanate from the internal path instead of arc $(v_{i_{1, x}^-}, v_{i_{1, x}^+})$.

Case (c) : If only one expanded node in E_i , say $v_{i_{1, x}^+}$ is visited by T'' , then select $v_{i_{1, x}^+}$ as the first node and another arbitrary expanded node of v_i as the last node, say $v_{i_{q, x}^-}$. Take the arc $(v_{i_{1, x}^+}, v_{i_{q, x}^-})$ to emanate from the internal path, and step (1.c) is skipped.

(1.c) Replace the arcs directed to and emanated from the paths $(v_{i_{1, x}^-, \dots, v_{i_{1, x}^+}), \dots, (v_{i_{q, x}^-, \dots, v_{i_{q+1}, x}^+})$ with the arcs $(v_{i_{1, x}^-, \dots, v_{i_{1, x}^+}), \dots, (v_{i_{q, x}^-, \dots, v_{i_{q+1}, x}^+)$, i.e., connect the first node and the last node of each internal path together.

(2) For each pair of mutual exclusive nodes v_i and v_j do followings.

If the covering cluster set of v_j is contained in the covering set of v_i , then bypass node v_j by connecting the preceding node to the succeeding node of v_j .

For example, the new g-tours of the g-tours in Figure 5.11 and Figure 5.13 are shown in Figure 5.14 and Figure 5.15 respectively. The cost of the resulting g-tour depends on the strategies in interrelated steps of algorithm CAN-TRANS. The best solution can only be made by comparing all possible alternatives. Ad hoc strategies are preferable to an optimal one since this procedure itself is a part of a heuristic algorithm.

Although the arcs of I-N graph does not satisfy the triangle inequality (which is discussed in next section), the canonical transformation will guarantee that the cost of the resulting g-tour will not be more than that of the original g-

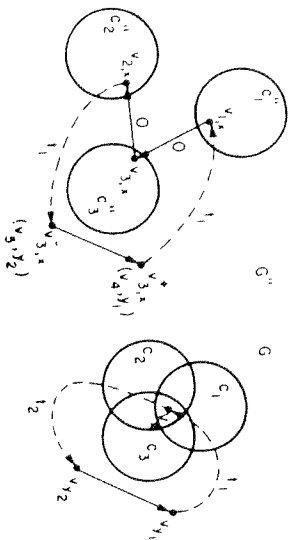


Figure 5.14 The new g -tour of T''' in Fig. 5.11 after canonical transformation.

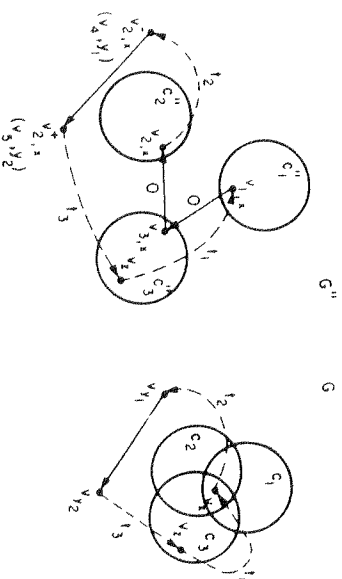


Figure 5.15 The new g -tour of T''' in Fig. 5.13 after canonical transformation.

tour as long as the arcs of the original graph satisfies the triangle inequality. This is proved in the next theorem.

Theorem 5.8. If the arcs of a graph G satisfy the triangle inequality, then the canonical transformation transforms any g -tour T''' on the L^N graph G''' into another g -tour T in G with equal or less cost.

Proof: Since arcs in G satisfy the triangle inequality, step (2) of Algorithm CAN-TRANS will transform a g -tour into another g -tour with no more cost. We only need to show that steps (1.b) and (1.c) transform a g -tour on G''' into another with no more cost (step (1.a) will not change T'''). Suppose v_i is the node under consideration, and the expanded nodes of v_i that are contained in T''' are connected as the following internal paths, $(v_{1,x}, \dots, v_{q,x})$, $(v_{1,x}, \dots, v_{q,x})$, \dots , $(v_{p,x}, \dots, v_{q+1,x})$.

Case (a), (b), and (c) in the following refer to Case (a), (b), and (c) respectively in step (1.b) of Algorithm CAN-TRANS. First, if they are in Case (a), then the arcs $(v_{1,x}^-, v_{1,x}^+)$ and $(v_{1,x}^-, v_{1,x}^+)$ are kept unchanged in the new T''' . However, the arcs $(v_{1,x}^-, v_{1,x}^+)$ and $(v_{1,x}^-, v_{1,x}^+)$ are replaced by $(v_{1,x}^-, v_{1,x}^+)$, the arcs $(v_{1,x}^-, v_{1,x}^+)$ and $(v_{1,x}^-, v_{1,x}^+)$ are replaced by $(v_{1,x}^-, v_{1,x}^+)$ and so on. The replaced arcs are all existing arcs in G , so they should satisfy the triangle inequality. Each replaced arc is either an existing arc in G or a new one with zero cost. The replacement will not increase the cost of T''' in both cases. Further, the cost of any internal tour is zero. Therefore, the cost of the new g -tour is not increased. Second, if they are in Case (b), then the arc $(v_{1,x}^-, v_{1,x}^+)$ is kept unchanged in the new T''' and the cost of $(v_{1,x}^-, v_{1,x}^+)$ is

the same as the cost of $(v_{i,x}, v_{i,z}^-)$. The replacement of other arcs will not increase the cost for the same argument as in Case (a). So, the cost of the new g-tour is not increased. Third, if they are in Case (c), then the arc $(v_{i,x}, v_{i,z})$ is kept unchanged in T' , and the arc $(v_{i,x}, v_{i,z}^-)$ is replaced by the arc $(v_{i,x}, v_{i,z}')$, which has the same cost as that of $(v_{i,x}, v_{i,z}^-)$. Therefore, the cost of the new g-tour is the same as that of the old one. \square

In summary, a transformation consisting of two transformational phases transforms the GTSP on a graph with intersecting clusters into the TSP of another graph. The optimal GTSP solution to the original graph can be obtained from the optimal TSP solution to the transformed graph. The canonic transformation transforms any TSP tour on the transformed graph back to a canonic g-tour on the original graph, such that the heuristic TSP algorithms solves the GTSP.

5.6 Solutions For The Transformed Graph

In this section, we examine the properties of the transformed graphs (G-S graph) and suitable heuristic algorithms for solving TSP on these graphs. The properties of a G-S graph are as follows.

- (1) The graph is asymmetric.
- (2) The arcs do not satisfy the triangle inequality.
- (3) The graph is not a complete graph.

The arcs do not satisfy the triangle inequality since there exists arcs with zero cost. The graph is not complete since certain arcs are omitted to force the TSP

tour to be a feasible and canonical g-tour on the original graph. However, the graph can be transformed into a complete graph by adding dummy arcs with a very large cost. Heuristic TSP algorithms applicable only to complete graphs can then be applied to the resulting graphs.

Different approaches to solve the TSP optimally have been studied extensively in the literature [BEL68] [BCHR79] [CUP68]. These include (a) integer programming [DAN54] [MIL76] [MIL78], (b) dynamic programming [HEI62] [HEI62] [GON62], and (c) branch-and-bound [LIT63]. However, since the TSP is NP-hard, it is difficult to apply these algorithms to solve the GTSP optimally. For this reason, it is more realistic to apply polynomial-time heuristics for the TSP to solve the GTSP.

Over the past fifteen years, a large number of polynomial-time heuristics have been proposed to solve the TSP. These heuristics do not guarantee an optimal solution but seem likely to lead to "good solutions." These heuristics can be classified into the following categories.

- (1) Tour building : Starting from any node, a tour is built by successively including other nodes into the tour. *Nearest neighbor* [BEL68], and *insertion* (arbitrary insertion, nearest insertion, cheapest insertion, farthest insertion) [KAR64] are two examples.
- (2) Tour-to-tour improvement : Starting from a tour, better tours are obtained by changing the order of some nodes iteratively. An example is the *K-exchange* [LIN65] [CHR72].
- (3) Partitioning and decomposition: A large problem is decomposed into a series of subproblems first, then each subproblem is solved independently. *Strip* is such an example [BEA59] [BEN83].

- (4) Relaxation: The TSP is relaxed into another problem, such as the minimum spanning tree [HEL70] [HEL71] [HON76], optimal assignment problem [KAR77] [KAR79] [GIL64], n -path [HOL80], and matching problem [BEL71], then the solution is transformed into a TSP tour.

- (5) Others: This category includes the longest-path-problem transformation [HAR62] [PAN64], analog computer approach [CLA78], etc.

Given a heuristic H , let $R_n(H)$ be the ratio of the obtained tour length to the minimal tour length. For the TSP unconstrained by the triangle inequality, for any $k > 1$, the problem of finding a tour with R_n of length bounded by k is NP-complete [GAR79]. As a result, it is unlikely that polynomial-time heuristic with constant upper bound on $R_n(H)$ exists. However, for the asymmetric TSP unconstrained by the triangle inequality, Karp [KAR79] gave a polynomial-time heuristics with complexity $O(n^3)$. If the distances of the graph are drawn independently from a uniform distribution, with probability tending to 1, this heuristic has $R_n < 1 + \epsilon(n)$, where $\epsilon(n)$ goes to zero as $n \rightarrow \infty$.

Although the distance of the transformed graph of the GTSP do not satisfy the above assumption, the algorithm still work well for the transformed graph. We have evaluated this algorithm on a set of test problem instances generated by randomly clustering cities in a map of the Unite States, which is the same as the 33-city problem in Karp's paper [KAR64]. In Karp's algorithm [KAR80], the problem is first solved as an assignment problem. The resulting assignment may contain some cycles. A patching algorithm is then applied on these cycles to produce a feasible TSP tour. The patching algorithm iteratively selects two cycles and replaces an arc from each cycle with a pair of arcs that connect the two cycles together. However, since the transformed graph is not complete, the patching may not always be possible. Dummy arcs with a very large cost may be

included into the patched cycles. These arcs are likely to be replaced by other arcs in later iterations. However, if the result of patching still contains such dummy arcs, then the resulting g -tour is permuted by an exhaustive search algorithm to eliminate the dummy arcs. Figure 5.16, 5.17, 5.18, 5.19 and 5.20 are the results on some test graphs whose cities are clustered in a more or less geographically related way. Table 5.1 shows the results of a set of test graphs whose cities are clustered randomly. The number of clusters varies from 3 to 8. The average values shown in the table are the average on the ratios of the heuristic solution to the optimal solution of all instances. The optimal solution of a problem instance is obtained by a branch-and-bound algorithm. Due to the lack of a good bounding estimate for graphs with intersecting clusters, the clusters in all testing graph are non-intersecting. On the VAX 11/780 computer, in most cases, it took about nine hours CPU time on the average to get an optimal solution for a graph with eight non-intersecting clusters. The long execution time is due to the inability of purging large number of solutions in the branch-and-bound algorithm.

The following example is to show how the query processing problem in Figure 5.1 and 5.2 is solved. The graph representation of the problem is shown in Figure 5.21. In the problem, a query initiated from site X is to join R_1, R_2, R_3, R_4 together. These files are allocated at sites W, X, Y, Z . In the graph, each site is a node; the cost of arcs are the communication cost between two sites; and the sites storing the same file are grouped into the same cluster. For simplicity, the graph is assumed symmetric and arcs are bidirectional with the same communication cost in both directions. As discussed in Section 5.1, the query processing problem can be solve as a GTSP problem if the communication

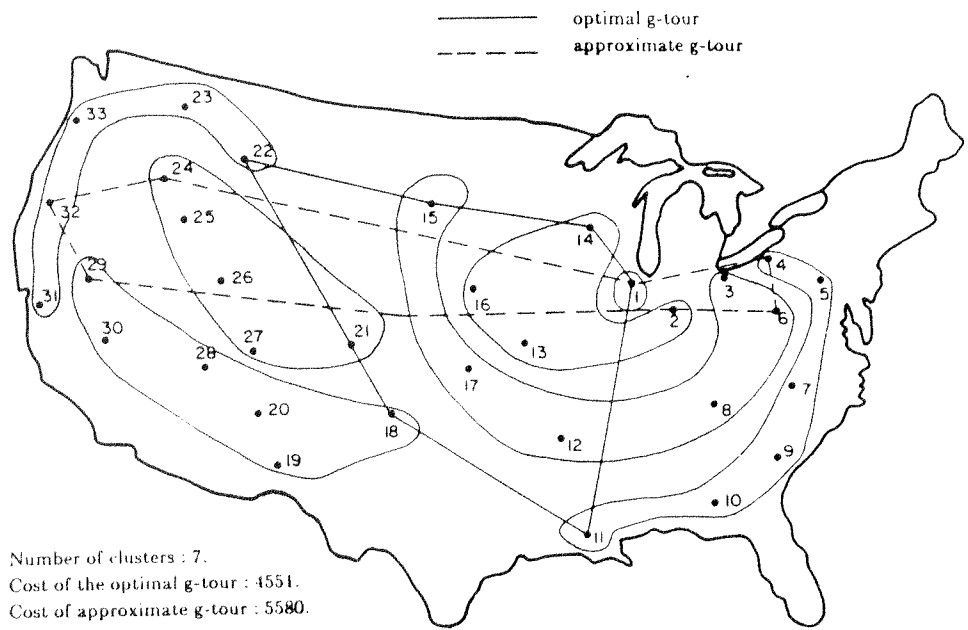


Figure 5.16 Example graph 1: optimal and approximate GTSP solutions.

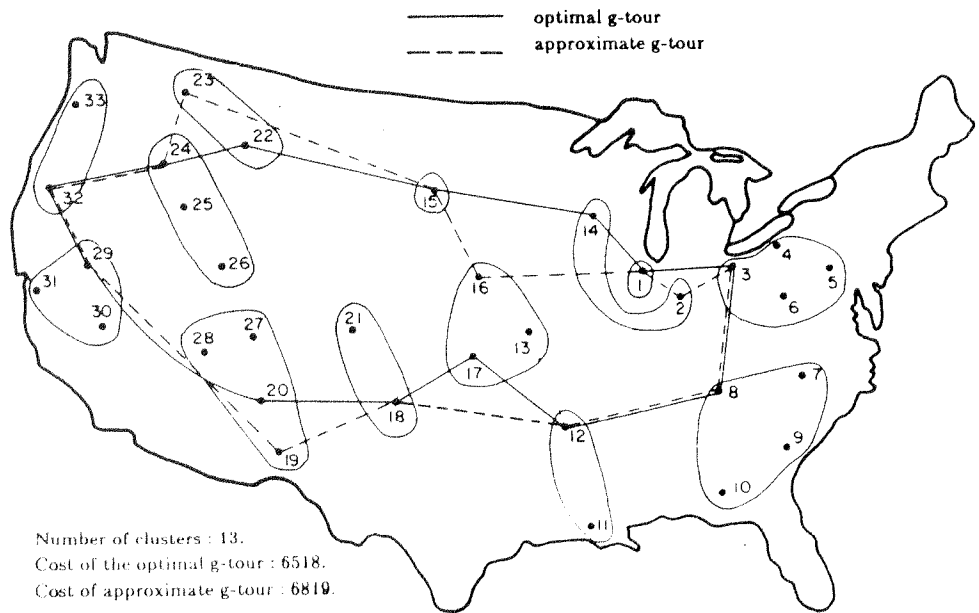


Figure 5.17 Example graph 2: optimal and approximate GTSP solutions.

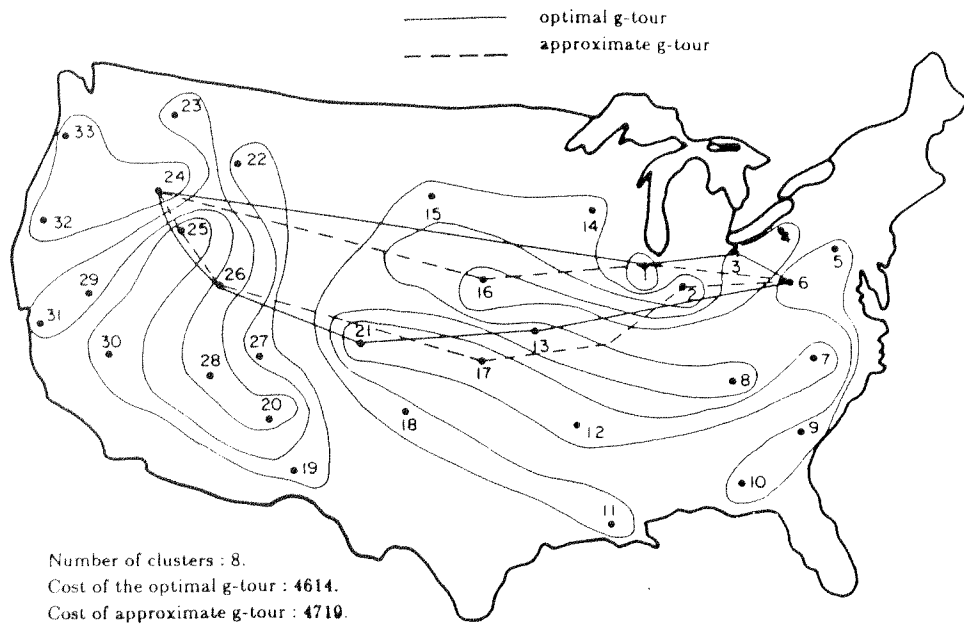


Figure 5.18 Example graph 3: optimal and approximate GTSP solutions.

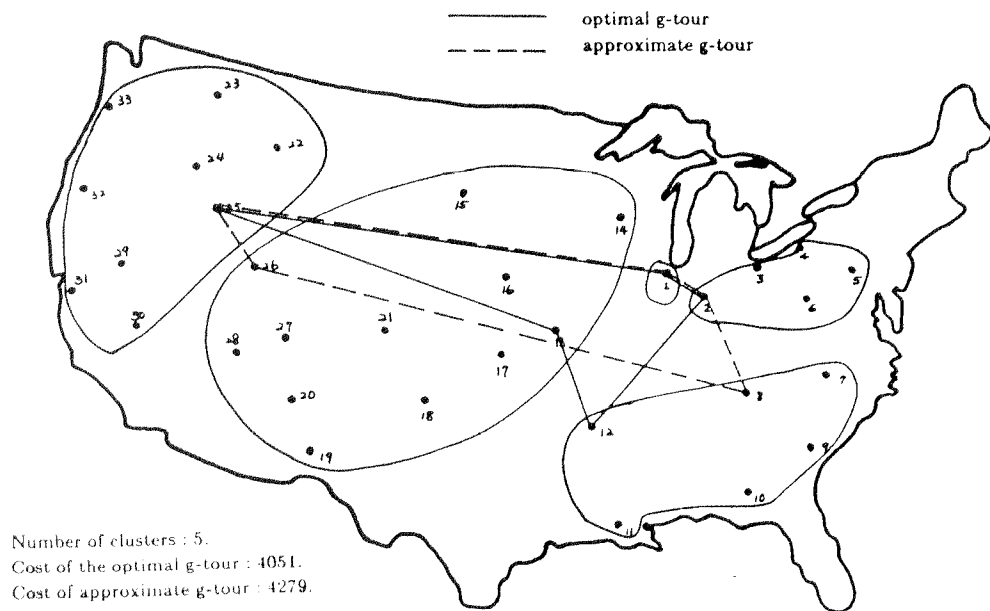


Figure 5.19 Example graph 4: optimal and approximate GTSP solutions.

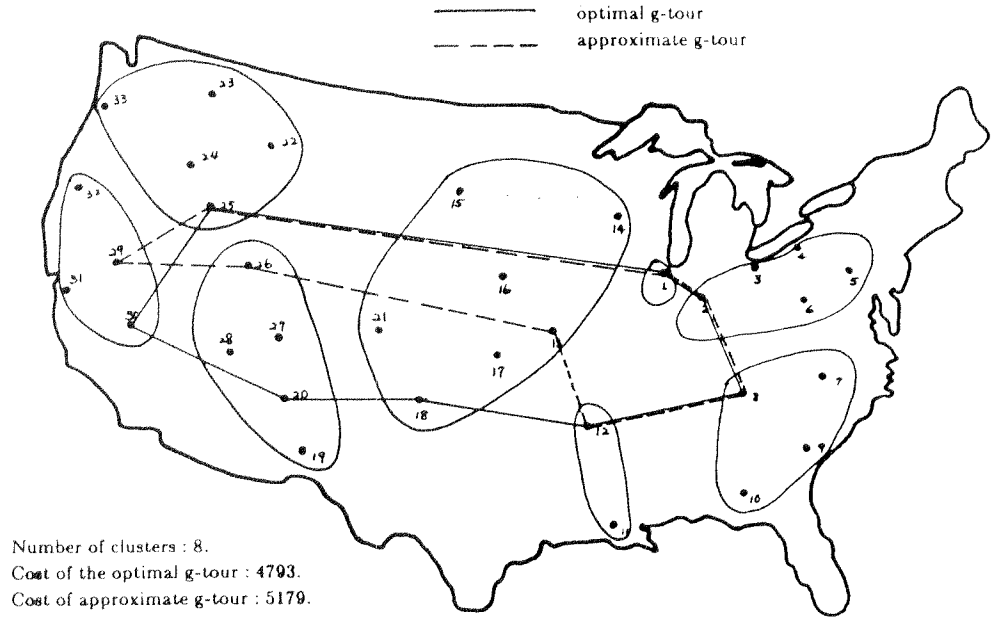


Figure 5.20 Example graph 5: optimal and approximate GTSP solutions.

Table 5.1 Performance of the TSP-transformation approach as solved by Karp's TSP algorithm.

Number	Cost of G-tours			
	Average	Standard Dev.	Maximum	Minimum
3	1.00163	0.001551	1.01635	1
4	2.35578	0.331121	3.54386	1
5	1.58284	0.153909	2.37311	1
6	1.95684	0.133981	2.62162	1.1785
7	1.52378	0.117846	2.16867	1.03278
8	1.40966	0.109277	2.00763	1.01676

• Optimal cost is 1
 • 10 cases are executed for each type of graphs with a certain number of clusters

Number of clusters	Execution Time									
	Average Karp	Average Opt.	Stand. Dev. Karp	Stand. Dev. Opt.	Maximum Karp	Maximum Opt.	Minimum Karp	Minimum Opt.	Maximum Karp	Maximum Opt.
3	16907	7	842	0.7	23778	12	1478	4.5	19885	36
4	16310	21	319	2.2	18085	36	15003	13.7	17121	208
5	16102	83	243	16.7	17306	67.4	14907	15.1	17306	67.4
6	16328	306	140	55	17306	2083	15892	112	17306	2083
7	25343	4888	1605	2083	33147	18392	17594	280	33147	18392
8	24034	32218	2795	15980	42595	159997	16401	849	42595	159997

• CPU time unit : seconds on a VAX-11/780

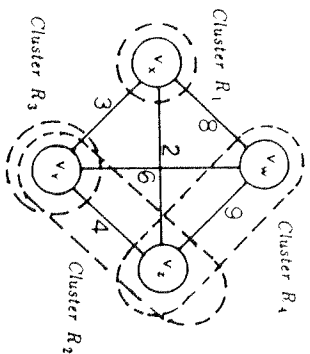


Figure 5.21 Graph representation of Figure 5.2.

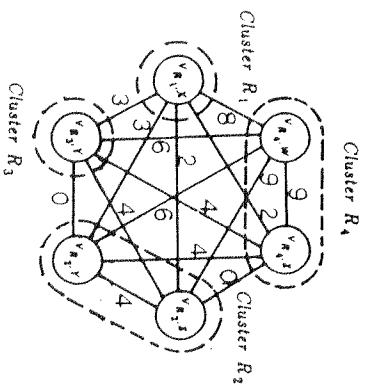


Figure 5.22 I-N graph of the graph in Figure 5.21.

cost between each pair of nodes is independent of the volume of transferred data. The corresponding I-N graph and G-S graph are shown in Figure 5.22 and 5.23 respectively. For simplicity, single-node-clusters, clusters R_1 and R_3 , are not expanded in Figure 5.23, and the arcs connecting expanded nodes are bidirectional. One of the TSP-tour in the G-S graph is shown in Figure 5.24. The corresponding tour and g-tour in the I-N graph and the original graph are shown in Figure 5.25 and 5.26 respectively. The joining sequence shown in Figure 5.1 can then be obtained from Figure 5.26.

5.7 Summary

In this chapter, we have examined the characteristic of communication networks in the future and have proposed an algorithm to solve the distributed query processing problem on such networks. The problem is then formulated as a generalized traveling salesman problem. The properties of the generalized traveling salesman problem and various approaches to solve it are studied. A set of transformations have been proposed to solve GTSP by applying the solution algorithms of TSP on the transformed graph. A graph with intersecting clusters is first transformed into another graph with non-intersecting clusters. The only requirement for this transformation to be valid is that the triangle inequality must be satisfied in the original graph. The graph with non-intersecting clusters is then transformed into a third graph such that the solution of the standard traveling-salesman problem on the third graph provides a solution to the original problem. The transformation is done by expanding each node in the original graph into three nodes and adding m nodes to control the routing. For

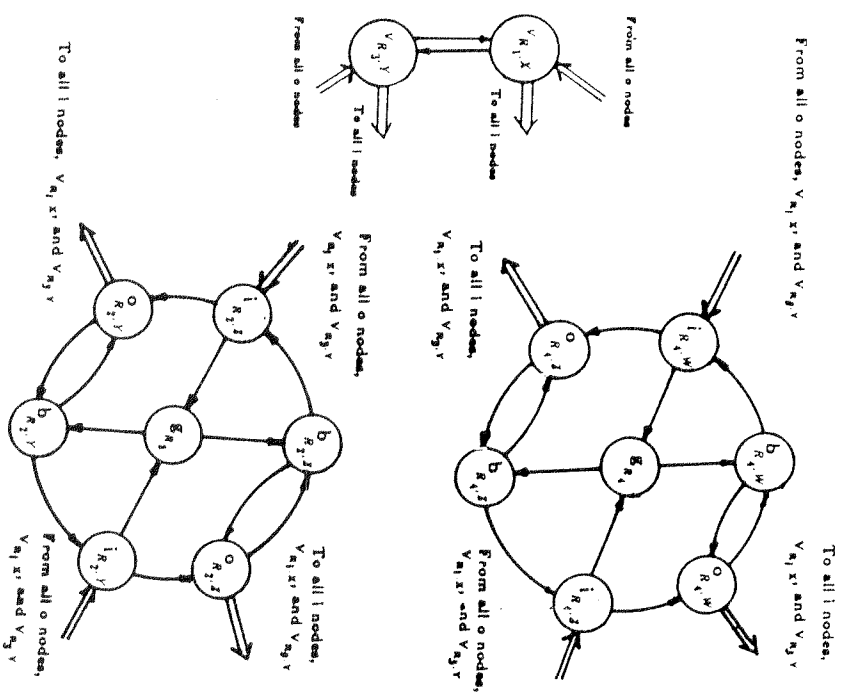


Figure 5.23 G-S graph of the graph in Figure 5.22.

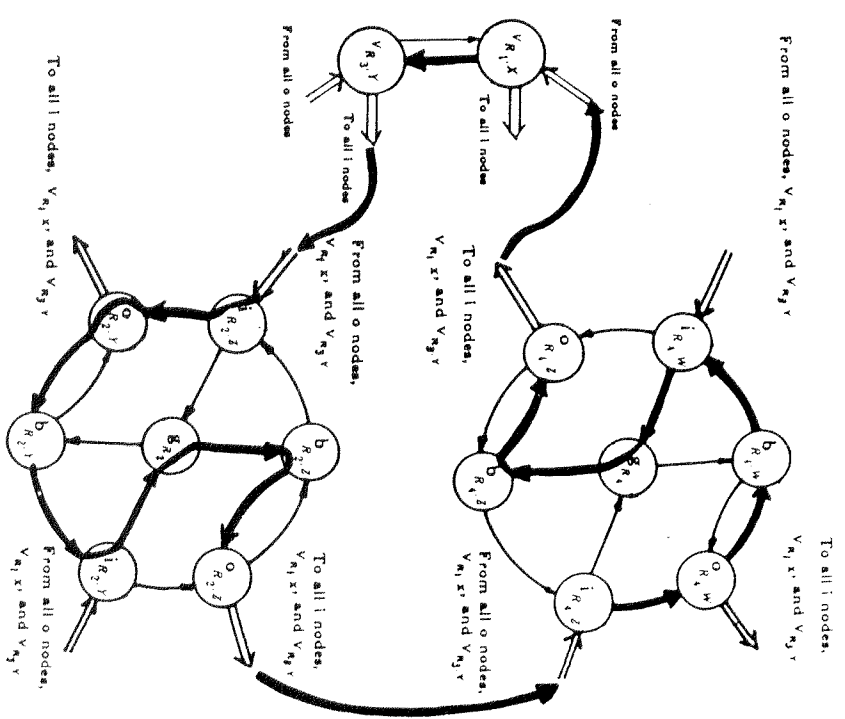


Figure 5.24 A TSP tour found from the G-S graph.

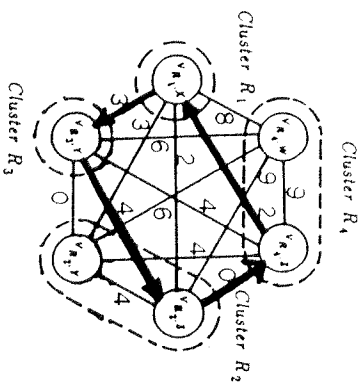


Figure 5.25 The TSP tour on the $I-N$ graph.

graphs with intersecting clusters, another transformation converts the possibly infeasible solution obtained by solving TSP on the transformed graph into a feasible and canonic solution of the original graph.

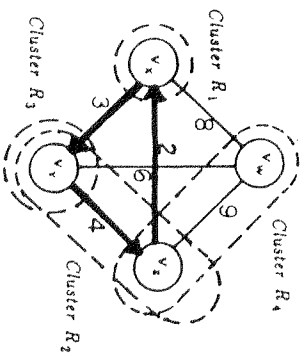


Figure 5.26 Final G-tour on the original graph.

CHAPTER VI

PERFORMANCE EVALUATION

One of the best ways to evaluate the DDHBMN is to implement a prototype system and then benchmark its performance. However, such a study would need extensive time and resources, which is far beyond the limit of this thesis. Further, it's inefficient to adjust the design by building new prototypes. Instead, a simulation study is conducted to demonstrate the feasibility, correctness, and effectiveness of the DDHBMN. The performance of different processing strategies are compared in the simulations.

6.1 Objectives of the Simulation

The objectives of the simulations are presented in this section.

1. **Showing the Correctness.** The design of a complicated integrated system such as DDHBMN may not be bug-free even if each individual algorithm is proved correct. Further, the correctness of some algorithms in this thesis, such as cc_minimunlock, is hard to prove formally. Hence, simulations are necessary to demonstrate their correctness.
2. **Discovering the implementation difficulties.** Some assumptions made in this research may not be easy to implement in a particular system. For example, no restriction is made on the number of Transaction Handlers that a system can initiate concurrently. This, in turn, makes an implicit assumption that each

local system can initiate as many TH's as possible. However, this is very difficult to implement in a small computer, since the number of concurrent processes that a system can handle is usually limited. As a consequence, a transaction may be blocked after the initial broadcast in the home site due to the lack of available processes in other sites. Such problems can be easily discovered through simulations and then solved before the system is actually implemented.

3. Comparing different strategies. When selecting a control strategy among possible alternatives in implementing the DDBLMIN, the system designer should be aware of the differences among various strategies on the potential of achieving system goals under different conditions. Henceforth, a comparative study is one of the objectives in this simulation study. We have compared between

- (a) redundant and non-redundant materialization, and
- (b) releasing locks after the LP phase and after the RT phase.

These strategies are compared under the following conditions:

- (a) the ratio of the number of update transactions to the total number of transactions,
- (b) the number of transactions a site can handle,
- (c) the ratio of packet size in the communication system to the block size in disks, and
- (d) the transaction arrival rate.

6.2 Queuing Network Model

Figure 6.1 depicts the *event-driven* queuing model of DDBLMIN. The physical resources are modeled as servers, which include DISK's, CPU's and the NETWORK. The logical resources (relations) are not shown in Figure 6.1, which are controlled by the concurrency control protocols and cannot be modeled as a simple server. However, a pseudo server "THINK" is used to model the lumped waiting time in the transaction processing. A request for a service, the start of a service, or the termination of a service is an event. The routing of tokens is dependent on the state of the system. An active TH in a site is represented as a token, which circulates around the servers in a site and the server NETWORK. A token in a site is generated for an incoming transaction, which in turn generates another set of tokens in other sites after the initial broadcast. Each such token represents a TH in a related site. A token is removed either after the transaction terminates or after the site where it is located is no longer involved in the processing. A token is *busy* if it is being served by a physical server (a physical server represents a physical resource) *waiting* if it is waiting for service from a physical server; and *sleeping* otherwise. A token is also called *active* if it is not sleeping.

All processes in one site share the single CPU in a 'process sharing' mode. The service time of a CPU server is proportional to the complexity of the job and the volume of the data to be processed in the CPU. Since the processing time of the local database in each site is dominated by the disk accessing time, it is modeled as a DISK server working in FIFS (first come first serve) mode. Only one token can be serviced at a time, and other tokens requesting a DISK service should wait in the DISK queue. The service time of a DISK server is a function of the disk seek and data transfer times. The disk seek time is a

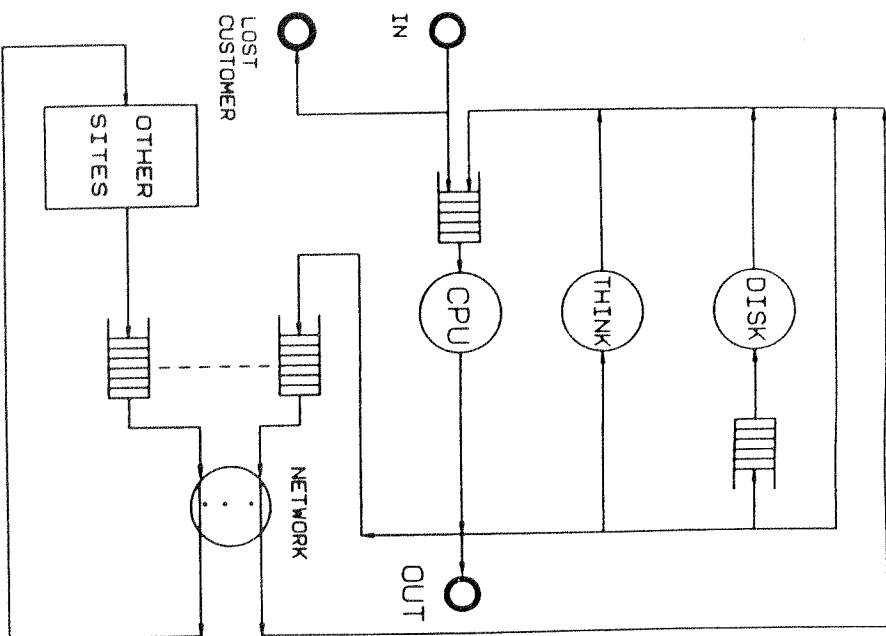


Figure 6.1 Queuing model of simulator DDBLMN-SIM.

uniformly distributed random variable. The data transfer time is proportional to the volume of data to be read from or written to the disk. All sleeping tokens in a site are put in the local server THINK until wakened up by other tokens. There is only one NETWORK server in the model, which is operated in a "random selection" mode to model the behavior of the CSMA/CD protocol. In each site, a first-in-first-out (FIFO) queue is used to store all local tokens (tokens generated from this site) that request NETWORK services. The NETWORK server randomly picks one of the non-empty queues and serves the first-arrived token in this queue. This type of servers is called a RANDOM server in this thesis. The service time is also proportional to the volume of data to be transferred.

6.3 Process Flow of DDBLMN

The process flow of a transaction in the system is shown in Figure 6.2. A transaction is initialized in the home site first. After the initial broadcast, all sub-queries are processed by a set of sites cooperatively. In each sub-query processing, a TH requests read-locks from the local Concurrency-Control subsystem for those relations stored locally and used in the LP phase. The Concurrency-Control subsystem grants the locks to the TH temporarily if the request passes the consistency checking. At the end of LP processing, a synchronization phase is proceeded to synchronize all sites. The synchronization process also serves as lock confirmation. After the synchronization process starts successfully, the temporarily granted locks are permanently granted to the transaction. Before the locks are confirmed, any conflicting lock broadcasts on the network will invalidate these temporary locks. The LP phase should be restarted in this case.

TRANS. PRO.

SUBQUERY PRO.

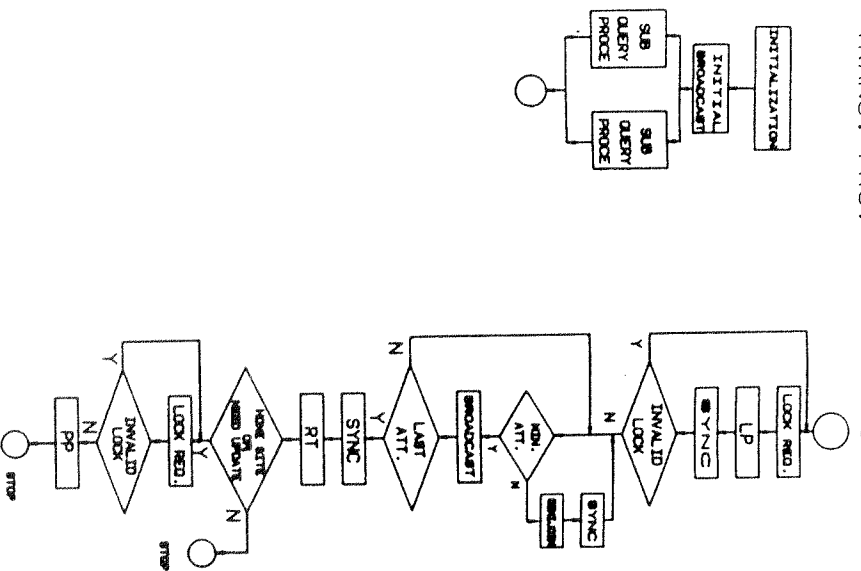


Figure 8.2 Process-flow diagram of DDBLMN-SIM.

Following the successful synchronization, the CSI phase starts immediately. During this phase, a sequence of global semi-joins are processed iteratively. In starting an iteration, a TH requests the read-locks for those relations not locked in the LP phase. The minimum attribute is then identified and broadcast to other sites which execute the semi-joins with the minimum attribute. The iteration then terminates, and the next iteration begins. The RT phase is started after all attributes are broadcast. All related sites broadcast in turn their relations that have not been broadcast by other sites to the home site.

A TH not in the home site terminates at the end of the RT phase if there is no relation in this site to be updated. All sites with relations to be updated continue their processing in the PP phase. The write locks are first obtained in this phase. After the write-locks are temporarily granted by the Concurrency-Control subsystem, the home site broadcasts the identity of the target data with new values to all sites. The write-locks are confirmed by the broadcast. All sites will update their database and terminate the transaction. The write-locks may be invalidated before the broadcast if there is a conflicting lock broadcast during this phase. The PP phase should be restarted in this case. All relations to be updated are still read-locked in this phase, therefore, there is no need to restart from the beginning.

6.4 Assumptions in The Simulation

To simplify the simulation so it can be conducted within a reasonable amount of time, the following assumptions are made.

- (1) There is only one query in each transaction.
- (2) There is at most one attribute to be updated.

- (3) The attribute to be updated is also a queried attribute. Hence, a write is always preceded by a read.
- (4) If non-redundant materialization is used, a copy of each queried relation is randomly selected. (In the simulations, the non-redundant materialization is really implemented by a non-replicated file allocation.)
- (5) After semi-joined of joining attribute of selectivity p , a relation (say, R_i) is reduced from the original size of s_i to a size distributed within 30% deviation of $s_i \times p$. The estimated size of a relation after a semi-join is $s_i \times p$. The indirect semi-join effect is ignored.
- (6) The heuristic function used in the GSI phase is assumed to be the product of the effective size and the relative selectivity of the candidate minimum attributes.
- (7) The join restriction effect in the RT phase is ignored since it is insignificant and is expensive to calculate in the simulations.
- (8) The transmission order in the RT phase is governed by the site numbers.
- (9) Since only one message can be broadcast from any site at any time, we assume that the Network-Interface subsystem in each site can only process one network access request at any time.
- (10) The service time of a CPU server is independent of the number of processes sharing the CPU. We make this assumption for two reasons. First, the CPU time is about one to two orders of magnitude less than the disks and network service times (the minimum overhead of the disk/network is a block/packet). Second, inexpensive multiprocessors are going to be available soon such that each TH may be served by a dedicated processor in such systems. (Intel has already produced a hypercube computer system

consisting of 32 80286 microprocessors with a cost less than \$100,000.) Finally, switching the CPU to service another token during the simulations is very computationally expensive.

6.5 Simulator DDBLMN-SIM

The simulation is implemented by a simulator DDBLMN-SIM, which is written in the C programming language and developed on a *ENVY** systems running *BSD 4.2** at the Department of Computer and Information Science, the Ohio State University. The simulator is executed on a SUN workstation running the same operating system.

A set of primitive subroutines in a simulation management package DDBLMN-SMPL, which is a modification of SMPL [MAC80] is used to control the simulations. These subroutines include initialization operations, event scheduling, facility reservations, random variable generation, data collection and reporting routines. A simulation model using DDBLMN-SMPL subroutines is a static network of FCFS and RANDOM facilities. Each of these facilities is associated with a queue in which incoming tokens wait for service. When the facility becomes available, a token is taken off the queue and serviced. The service involves reserving the facility for a given time period and, when that time has expired, passing the token to either another facility or the THINK state to sleep. After initialization of facilities, the simulator performs four basic operations iteratively.

- (1) Reserves a facility for a token to disallow any other users from using it. If the facility is already in use, the requesting token is placed in a queue waiting for the facility.

* UNIX is the trademark of AT&T. BSD 4.2 is Berkeley version 4.2.

- (2) Releases a facility, and allows a token (either the first one or randomly selected) in its queue to access to it if the queue is not empty.
- (3) Schedules a new event to happen at some future time. This action places the new event in a time-ordered queue of future actions. This step is referred to *new-event-scheduling*.
- (4) Causes the next event to occur. This action removes an event from the time-ordered queue.

Figure 6.3 shows the fundamental flow of events used in the simulations. After the queuing network is initiated, events are repeatedly generated until the following termination conditions are satisfied:

- (1) the system response time is stabilized, and
- (2) at least a certain number of transactions are executed.

Condition (2) is needed to increase the confidence on the property of deadlock free.

Each event asks the simulator to either reserve or release a facility. If the facility is successfully reserved, the simulator schedules a release event to occur at the end of the service time. If the event was to release a facility, then the program releases it, routes the token to the next facility according to the state of the system, and schedules a new event to reserve the new facility immediately.

A token in any queue can be removed and rerouted to the local THINK server. This is useful for modeling the behavior of a Transaction Handler when the temporary locks it holds are invalidated by another transaction.

It is not very difficult to implement the queuing network model in Figure 6.1 by DBL:MN:SMPL. For DISK servers, the implementation is

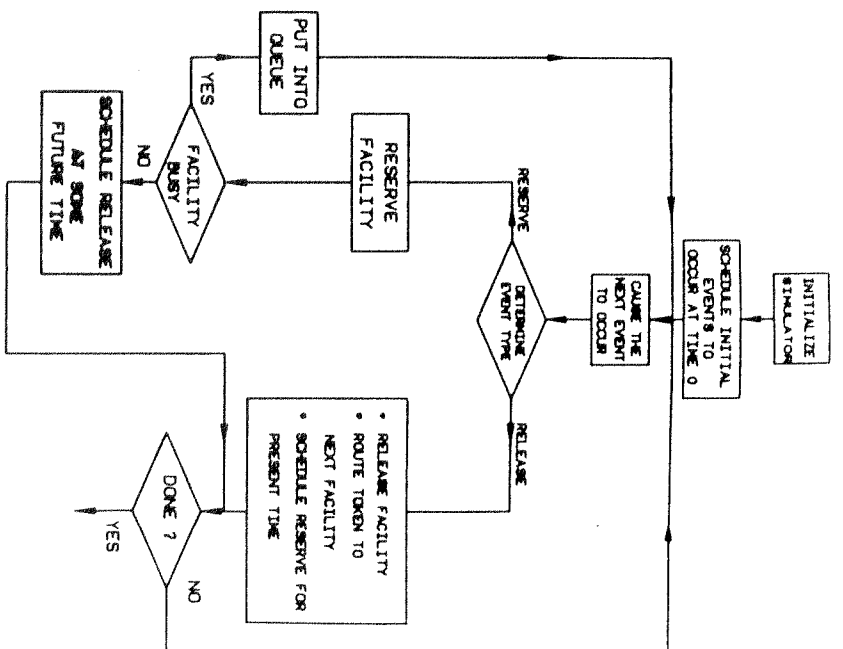


Figure 6.3 Event flow of DBL:MN:SMPL.

straightforward since they are FCFS servers. However, modifications to DDBLM:N-SMPL are needed to model resources operating in other types of service and queuing disciplines. For the servers that can service all incoming token parallelly, such as THINK and CPU, there is no associated waiting queue. The finish time can be immediately calculated. Thus, next event can be immediately scheduled when a service is requested. When a token needs a service from the CPU at time t_1 , the service time t_2 is first computed and added to t_1 . The next event after the CPU service is scheduled at time $t_1 + t_2$. When a token is routed to a THINK server (either blocked or halted) in new-event-scheduling, it stays there until wakened up by another token. There is no new event to be scheduled in new-event-scheduling in this case. For the NETWORK server, which serves in a "random selection" discipline with multiple FIFO queues, we add a logical server in front of each queue to model the Network-Interface subsystem in each site. This is due to Assumption (9) in Section 6.4. Another queue is attached to the NETWORK sever itself. A token leaving the Network-Interface server is routed to the NETWORK queue. Also, at most one token from each site can be routed to the NETWORK queue. In the mean time, Network-Interface is not released to service another token unless the token in the NETWORK queue is served. When the NETWORK server is free, a token in the NETWORK queue is randomly picked to be served and the Network-Interface subsystem from which the token came from is free to service another token in its queue. The tokens in the Network-Interface servers of other sites should keep waiting for a NETWORK service. The DDBLM:N-SMPL implementation of the queuing model shown in Figure 6.1 is detailed in Figure 6.4. As we can see, there is no FIFO queue associated with servers THINK and CPU. The Network-Interface subsystem in each site is modeled as a FCFS

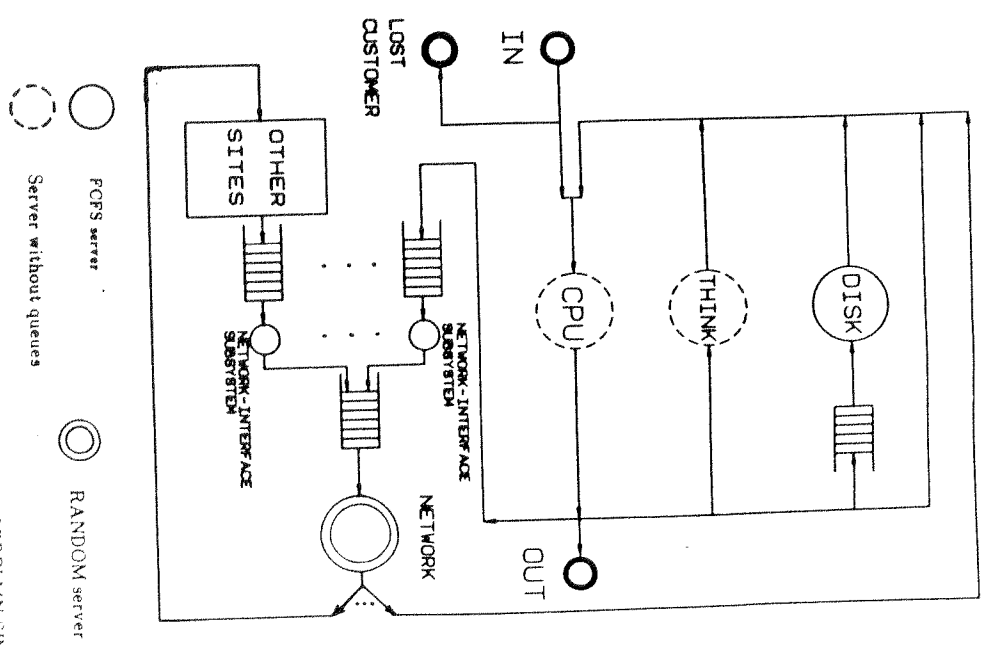


Figure 6.4 DDBLM:N-SMPL implementation of DDBLM:N-SIM.

Table 6.1 Parameters of DDBLMN-SIM system configuration.

Number of sites	4
MAX_TR_OF_SITE	5,8,10
Number of relations	4
Number of joining domains	4
Cardinalities of joining domains	80, 80, 160, 200
Width of joining domain	16, 8, 24, 40
Width of relations	240, 376, 408, 368
Initial cardinalities of relations	38707, 40960, 193536, 30720

MAX_TR_OF_SITE: max. number of transactions that can enter a site at the same time.

Table 6.2 Configuration of relations in DDBLMN-SIM.

(0/1: existence/nonexistence of an attribute in a relation).

	domain 1	domain 2	domain 3	domain 4
Relation 1	0	0	1	1
Relation 2	0	1	1	0
Relation 3	1	1	1	0
Relation 4	1	0	0	1

Table 6.3 Initial selectivities of all joining attributes in DDBLMN-SIM.

	attribute 1	attribute 2	attribute 3	attribute 4
Relation 1	-	-	0.6	0.9
Relation 2	-	0.4	0.9	-
Relation 3	0.8	0.5	0.1	-
Relation 4	0.8	-	-	0.5

Table 6.4 Relation allocation in DDBLMN-SIM.

	Site 1	Site 2	Site 3	Site 4
Relation 1	1	0	1	0
Relation 2	0	1	1	0
Relation 3	1	0	0	1
Relation 4	0	1	0	1

Table 6.5 Service time of servers in DDBLMN-SIM.

(time unit : 10 * second)

CPU	3.0 * number of bytes
DISK	uniform [0, 20000] + 1000 * number of blocks
BUS	1.0 * bits/packet * number of packets

Table 6.6 Other timing parameters in DDBLMN-SIM.

(time unit : 10 * second)

Bus contention slot	50
Preprocessing time	1000

Table 6.7 Data format in DDBLMN-SIM.

Number of bits/byte	8
Disk block length	1000 bytes
Packet length (bytes)	100, 500, 1000, 2000

Table 6.8 Parameters used for query generation in DDBLMN-SIM.

QUERY_REL_HIT	0.8
TARGET_REL_HIT	0.4
QUERY_ATT_HIT	0.4
TARGET_ATT_HIT	0.4
UPDATE_QUERY_HIT	0, 0.1, 0.7, 1
UPDATE_ATT_HIT	0.4

QUERY_REL_HIT: probability that a relation is queried by a query.
 TARGET_REL_HIT: probability that a relation is to be retrieved by a query.
 QUERY_ATT_HIT: probability that an attribute is queried by a query.
 TARGET_ATT_HIT: probability that an attribute is to be retrieved by a query.
 UPDATE_QUERY_HIT: probability that a query is an update.
 UPDATE_ATT_HIT: probability that an attribute in an update is to be updated.
 UPDATE_QUERY_HIT is the probability that a query is an update.
 UPDATE_ATT_HIT is the probability that an attribute in an update is to be updated.

All timing parameters are shown in Table 6.5 and 6.6. At a speed of 1 MIPS, we assume that a CPU can process a byte (8 bits) in 3 instruction cycles and that the processing time is proportional to the number of bytes to be processed. The disk access time consists of two parts: the block seek time and the data transfer time. The length of each block is 1000 bytes. The block seek time is a uniformly distributed random variable between 0 and 20 milliseconds, which represents the time of moving the disk head and locating the sector. It is assumed that there is no seek time between adjacency blocks transferred in the same request. Once the first block is located, successive blocks can be transferred continuously with a speed of 1 millisecond per block. The bus communication overhead is represented by a constant and a variable proportional to the volume of transferred data. We assume that all necessary initialization work is done in the Network-Interface subsystem, and the contention overhead is accounted for separately. Consequently, the service time of the BUS server is only dependent on the number of transferred packets. The number of overhead bits is assumed to be 208, which is the same as the Ethernet specification. A contention slot in the CSMA/CD protocol is assumed to be 50 bit-times. The preprocessing phase lasts for 1 millisecond. The data format is listed in Table 6.7. The size of a disk block is fixed to be 1000 bytes and 8 bits per byte, while the packet length varies from 100, 500, 1000, to 2000 bytes.

There is a single queue of transaction arrivals with an exponentially distributed interarrival time. The mean interarrival time (MIT) varies from 0.1, 0.3, 0.5, to 0.7 seconds. The incoming transactions are randomly routed to a site. A transaction to a site is lost if the number of pending transactions at this site has reached the maximum, which is a parameter defined to the simulator. The query pattern of a transaction is also randomly generated. Table 6.8 shows the

parameters controlling the query pattern generation. We call a relation/attribute as a *hit relation/attribute* if it is involved in the target data identification in a query. QUERY_REL_HIT/QUERY_ATT_HIT are the probabilities of a particular relation/attribute hit by a query. A hit relation must contain at least one hit attribute. The relation containing at least one hit attribute should be a hit relation. We call a relation containing at least one target attribute as a target relation. TARGET_REL_HIT/TARGET_ATT_HIT are the probabilities of a particular relation/attribute to be a target relation/attribute in a query. UPDATE_QUERY_HIT is the probability that a query is an update query. A queried relation is randomly picked to be updated if the query is an update. UPDATE_ATT_HIT is the probability that a target attribute is to be updated in a relation to be updated. In an updated relation, at least one attribute is to be updated. All these probabilities are set to 0.4 except that QUERY_REL_HIT is 0.8 and UPDATE_QUERY_HIT is varying from 0.0, 0.1, 0.7, to 1.0. All queries are read-only when UPDATE_QUERY_HIT is set to 0.0, and are updates when it is set to 1.0. Those tunable variables are summarized in Table 6.9. Each combination is run with a rates of mean interarrival time of 0.1, 0.3, 0.5, and 0.7 seconds.

6.7 Simulation Results

Due to resource contentions in distributed systems the performance measure of transaction response times is an autocorrelated quantity [FIS78]. Hence, the initial conditions of the simulation model may produce transient measurements in the simulation outputs. We ran a series of test simulations to determine the number of response time observations to be truncated in the beginning of each simulation run. By observing the resulting graphed data, we

Table 6.9 Tunable variables in DDBI.MN-SIM.

MAX_TR_OF_SITE	5,8,10
PACKET_SIZE (bytes)	100, 500, 1000, 2000
UPDATE_QUERY_HIT (or UQH)	0, 0.4, 0.7, 1
R_LOCK (or r_unlock_time)	0, 1
RED	0, 1

PACKET_SIZE: packet length in bytes

R_LOCK: the time that a read lock can be released

0: after RT phase

1: any time after use

RED: materialization strategy

0: non-redundant

1: redundant

chose the truncation point on the graph where it appeared that the mean response time measurement had achieved steady-state stability. For the data shown in Figure 6.6, the value 80 is found to be a good truncation point, which is the minimum number of transaction to be observed. A value 2% is a good allowed deviation, which is the difference between the simulated response time of current terminated transaction and the average simulated response time of all terminated transactions in the simulation run. Notice that there are some unstable cases, in which the systems are saturated, in Figure 6.6. The minimum number of transactions to be executed in a simulation is set to be 100 for the safety margin, and the allowed deviation is set to be 0.5%. We now present the simulation results of different combinations of parameters generated from Table 6.9. It must be emphasized that the numerical values found in this section are accurate only to the extent that the data values and assumptions used to design the simulator are representative of an actual system. The main concern of this simulation is to show the correctness of control algorithms, to discover the implementation difficulties, and to compare the relative performance of different control strategies as discussed at the beginning of this chapter.

8.7.1 Deadlock-free and Implementation Difficulty

Since it is impossible to prove the deadlock-free property of the system by simulation, demonstrating it in the simulation is complementary work to the proof in Chapter IV. During the testing stage of the simulator development, we have tried to push the simulator to the extreme cases to find any deadlock. The incoming rate of transactions is increased to an extent far beyond the capacity of the simulated system, the percentage of updating queries to the total number of queries is increased to 100%, and the simulator is run up to thousands of CPU

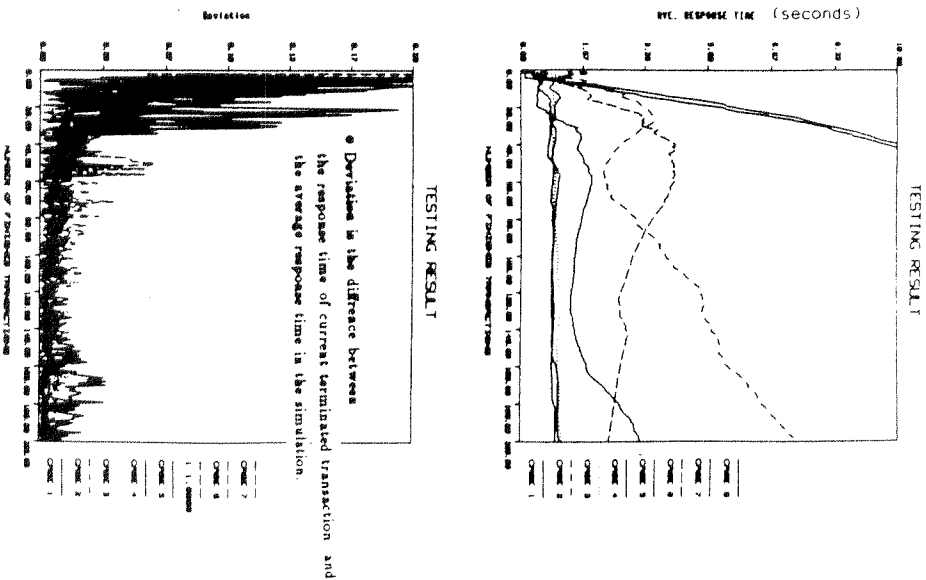


Figure 6.6 Testing results of DDBL/MN-SIM
 (a) mean response time; (b) deviation of response time.

minutes. There is no deadlock detected. We believe that these are sufficient to demonstrate the correctness of the concurrency control algorithms.

The implementation of DDRLMN is fairly complicated. We will not discuss the programming difficulty of the system since it is a common problem of any complicated working system. However, we would like to point out following problems.

(1) The simplicity of the concurrency control algorithms does contribute to the success of the simulator design. The major difficulties of the program development are on the state transition and the token routing, instead of on the concurrency control itself. Simplicity of the design is important in the implementation of a real system. We believe that DDRLMN would be a good choice to implement a DBMS in local environment.

(2) The system requires each site to keep track of the complete state of the database. To save communication overhead, state transition is really calculated by each site independently based on the network messages received. The complexity of such a distributed operation is fairly high. This not only causes some processing overhead, but also makes the system development and maintenance difficult. This is especially true in a larger system. The system designer should implement this mechanism with care. The possibility of errors in this part is high and may also lead to possible deadlocks. Comparing this to other portions of the system, this should deserve more attention. For example, a mistake on the minimum attribute identification or on the cardinality estimation may only affect the performance, but a mistake in this part may cause errors to the transaction processing. We would recommend that piggybacking the necessary

information in the broadcast message would help to reduce the complexity of the design and processing overhead. This is particularly true when the communication overhead is independent to the volume of data transfer.

(3) The complexity of computing the transitive closure is $O(N^3)$, which is a substantial overhead when the system is large. The most commonly used algorithm to compute the transitive closure is a matrix operation, in which the precedence graph is represented as an adjacency matrix [AHO74]. The overhead of this algorithm may be wasteful when the probability of conflict among transactions is low. For example, in an ATM system, customers do not issue conflicting transactions frequently. In this situation, the adjacency matrix representing the precedence graph is most likely sparse. Furthermore, since the transitive closure is always kept in the precedence graph, the changes to the precedence graph would not be substantial whenever a new precedence relationship is imposed or a node is removed. There is no need to recompute the precedence order for the entire graph. The algorithm discussed in Chapter IV can be modified as follows.

Algorithm DDRLMN-TRANSITIVE

Consider the graph as a family of isolated acyclic graphs, with one and only one directed edge between any pair of nodes in the graph. Node N_i represents the transaction TR_i , as stated in Chapter IV.

- (a) Remove the node and all the edges connected to it when the transaction that this node represents terminates.
- (b) Whenever a new precedence relationship is imposed on TR_i and TR_j , (say, TR_i precedes TR_j),
 - (b1) connect N_i to N_j ;

(b2) connect all preceding nodes of N_i to node N_j and all succeeding nodes of N_j .

We emphasize that this algorithm may not be better than the matrix approach when the probability of conflict is high and the adjacency matrix is not sparse.

6.7.2 Observations and Suggestions

The numerical simulation results are shown in Table 6.10. Each combination generated from Table 6.9 was simulated. The average, maximum, minimum, and standard deviation of all 192 combinations of different MIT on the following observations are listed.

(a) Mean Response time is the average response time of all transactions in a simulation.

(b) Mean queue length of NETWORK is

$$\left(\frac{\sum_{\text{all sites}} (\text{NETWORK} - \text{INTERFACE in the site})}{\text{number of sites}} \right) + \text{mean queue length of NETWORK}$$

in a simulation.

(c) Mean queue length of DISK is the average value of the mean DISK queue length in all sites, i.e.

$$\left(\frac{\sum_{\text{all sites}} (\text{mean queue length of DISK in the site})}{\text{number of sites}} \right)$$

(d) Utilization of NETWORK is the busy time of the NETWORK over the total time in a simulation

Table 6.10 Simulation results of DDH.MIN.SIM

Mean Response Time				
MIT	Ave	Max	Min	Stand. dev
0.1	7.053	27.900	0.190	0.5410
0.3	7.739	15.900	0.180	0.2995
0.5	0.833	5.190	0.170	0.0723
0.7	0.182	2.080	0.140	0.0256

Mean NETWORK time (total trans. time)				
MIT	Ave	Max	Min	Stand. dev
0.1	0.071	0.388	0.018	0.0018
0.3	0.198	0.371	0.011	0.0017
0.5	0.284	0.442	0.143	0.0084
0.7	0.238	0.485	0.150	0.0017

Utilization of NETWORK				
MIT	Ave	Max	Min	Stand. dev
0.1	0.820	1.000	0.280	0.0290
0.3	0.585	0.999	0.105	0.0247
0.5	0.407	0.949	0.055	0.0209
0.7	0.291	0.791	0.040	0.0180

Mean (CPU time / total trans. time)				
MIT	Ave	Max	Min	Stand. dev
0.1	0.055	0.173	0.001	0.0041
0.3	0.104	0.209	0.007	0.0051
0.5	0.124	0.220	0.028	0.0041
0.7	0.131	0.220	0.040	0.0040

Mean Utilization of DISK				
MIT	Ave	Max	Min	Stand. dev
0.1	0.259	0.453	0.084	0.0082
0.3	0.128	0.185	0.084	0.0014
0.5	0.084	0.104	0.084	0.0007
0.7	0.057	0.080	0.045	0.0005

Local Transactions (trans. / second)				
MIT	Ave	Max	Min	Stand. dev
0.1	0.0221	2.081	0.000	0.0306
0.3	0.000	0.000	0.000	0.0000
0.5	0.000	0.000	0.000	0.0000
0.7	0.000	0.000	0.000	0.0000

Mean queue length of NETWORK				
MIT	Ave	Max	Min	Stand. dev
0.1	4.033	8.331	0.093	0.1851
0.3	1.740	0.689	0.012	0.1580
0.5	0.808	3.483	0.003	0.0865
0.7	0.197	1.598	0.001	0.0201

Mean queue length of DISK				
MIT	Ave	Max	Min	Stand. dev
0.1	0.108	0.388	0.002	0.0025
0.3	0.019	0.087	0.002	0.0008
0.5	0.008	0.021	0.002	0.0003
0.7	0.002	0.009	0.001	0.0001

Total CPU Time (SVN)				
MIT	Ave	Max	Min	Stand. dev
0.1	938.593	1914.800	184.900	21.1398
0.3	732.788	1389.800	457.900	13.9473
0.5	883.789	944.400	453.000	7.8510
0.7	847.708	895.700	458.900	7.1943

Time Unit : second

(e) Mean utilization of DISK is the average utilization of DISK's over all sites in a simulation.

(f) Mean ratio of NETWORK service time to the total transaction time is the average of

$$\left[\frac{\text{total time serviced by the NETWORK}}{\text{total transaction time}} \right]$$

of each individual transaction in a simulation.

(g) Mean ratio of DISK service time to the total transaction time is the average value of

$$\left[\frac{\text{total time serviced by the DISK's}}{\text{total transaction time}} \cdot (\text{number of sites}) \right]$$

of each individual transaction in a simulation.

(h) Mean ratio of CPU service time to the total transaction time is the average value of

$$\left[\frac{\text{total time serviced by the CPU's}}{\text{total transaction time}} \cdot (\text{number of sites}) \right]$$

of each individual transaction in a simulation.

(i) Rate of lost transaction per second is the ratio of total number of lost transactions to the total simulation time in a simulation.

(j) Total CPU time is the total CPU time used in a simulation on SUN workstations.

The following are observations and suggestions we draw from the simulation results.

(1) From Figure 6.7, 6.8, and 6.9, it is clear that the system can handle the load when the MTT is longer than 0.5 seconds, but will be saturated when

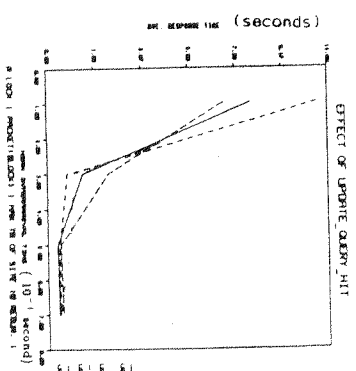
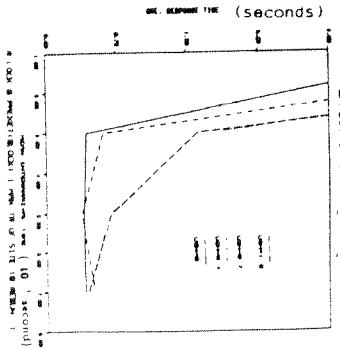
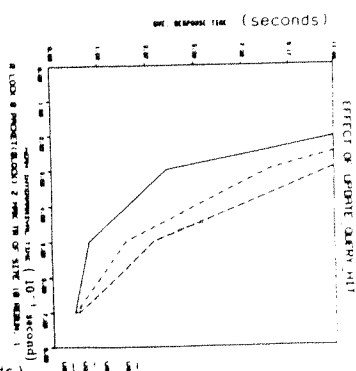


Figure 6.7 Mean response time.

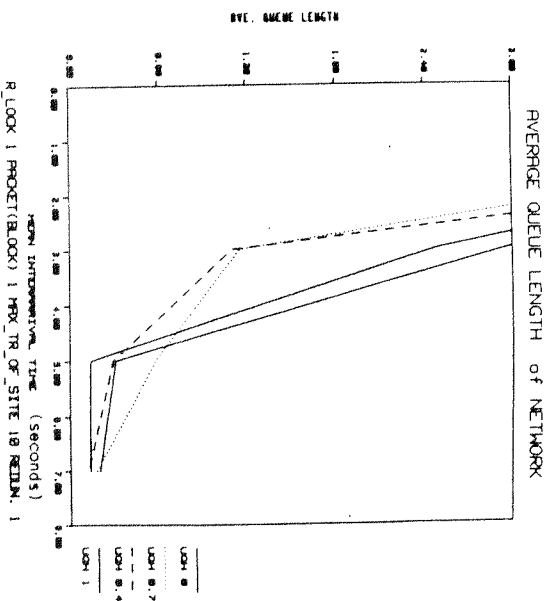


Figure 6.8 Mean NETWORK queue length.

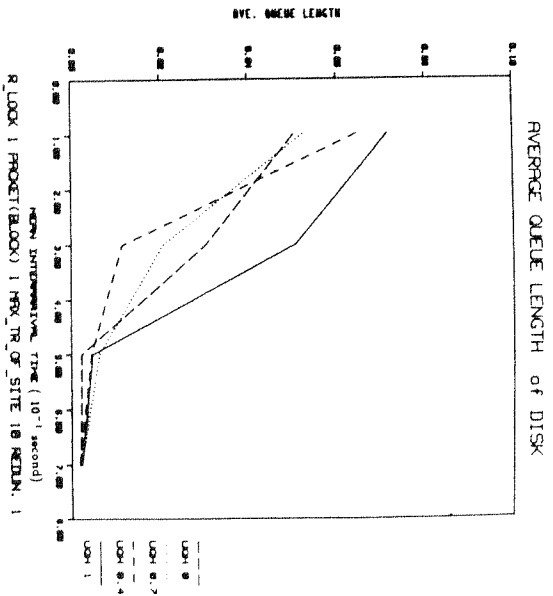


Figure 6.9 Mean DISK queue length.

the MITT reduced to 0.1 in many cases. When saturated, the queue of NETWORK (including NETWORK-INTERFACE) is full and a lot of incoming transactions are lost. This occurs even when all queries are not updates. When the MITT is longer than 0.3 seconds, the rate of lost transactions is almost invisible.

(2) Figure 6.7 shows that the percentage of updating queries does affect the response time. The more the updates are, the longer the response time.

(3) From Figures 6.8 and 6.9, we can also see that the UPDATE_QUERY_HIT has no significant effect on the length of either the DISK or NETWORK queues. That is true because we assume that at most one relation is updated and the volume of updated data is usually small. This fact also indicates that the communication overhead for lock management in this system is very low as we expected.

(4) Figure 6.10 shows the comparison of the utilization of DISK's and NETWORK. It is very clear that the NETWORK is heavily utilized, especially, when the system load is high. When the system is approaching saturation point, the utilization of the NETWORK is close to 100%. The DISK utilization is kept as low as 20% even at this point. This tells us that the NETWORK is the bottleneck of the system.

(5) Figure 6.11 shows the comparison of the average ratio of the simulated time serviced by CPU, DISK's and NETWORK to the total simulated transaction time. It is also very clear that a large portion of the transaction processing time is spent on the NETWORK. Only less than 10% of the time is spent on the DISK's and less than 5% on the CPU. Less than 50% of the time is waiting when MITT is kept longer than 0.3.

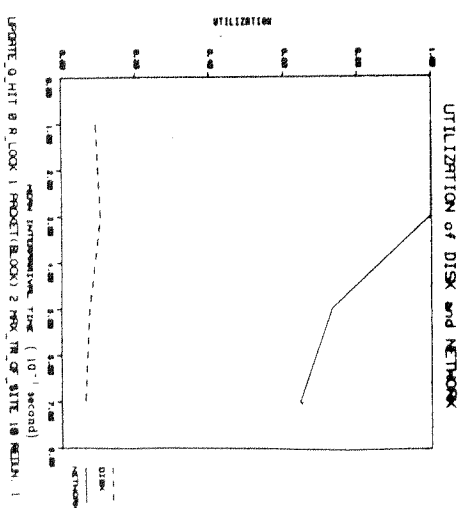
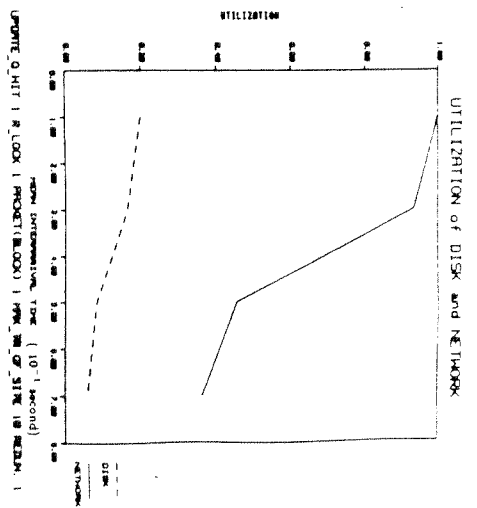


Figure 6.10 Comparison of the utilization of DISKs and NETWORK.

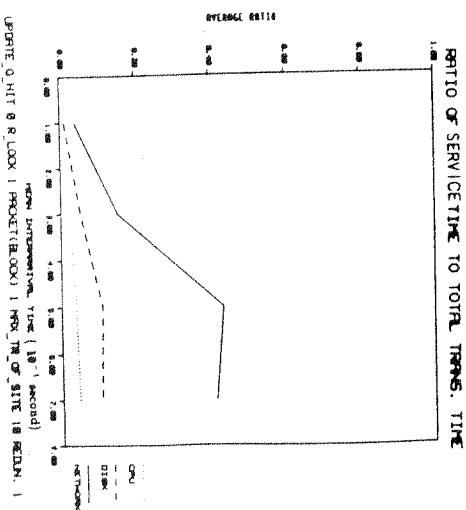
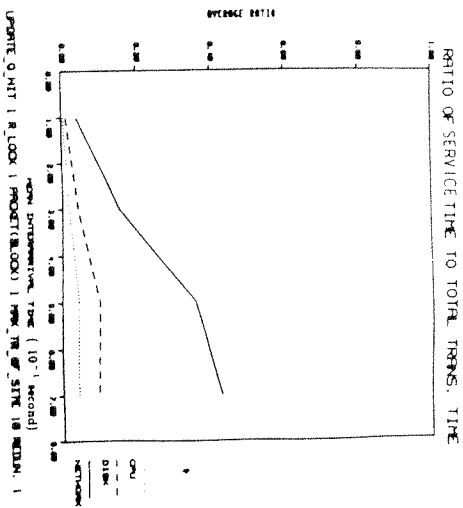


Figure 6.11 Comparison of the served time of CPU's, DISKS, and NET-
WORK.

(6) The size of packets in the NETWORK has significant effect on the response time as we can see from Figure 6.12. It seems unrealistic that the smaller the packets are, the better the system is. This is due to the fact that the fixed overhead in each packet is only 208 bits and the minimum packet length we simulated is 800 bits such that the fixed overhead does not affect the performance significantly. Really, there are many short broadcasts during the transaction processing. For fixed length packets, a large portion of their capacity is wasted.

There are two ways to reduce this overhead: either using variable length packets or piggybacking the necessary information into the sparse packets. Due to the time constraints, we did not simulate either case in our simulations.

(7) Figure 6.13 shows the effect of read lock releasing time. Two alternatives are simulated. When the UPDATE_QUERY_HIT is 0.0, in which all queries are read-only, all relations can be released right after use. Therefore, having a second copy in the working space has no advantage. Moreover, the induced overhead may even worsen the performance. As the UPDATE_QUERY_HIT is increased, having second copy does make the system performs better.

(8) The effect of materialization can be seen from Figure 6.14. Redundant materialization does have some improvement over non-redundant materialization. The difference is not great, this is because the scale of the simulation is too small such that at most two copies of a file are allocated. We believe that redundant materialization will be more attractive when the scale of the system is larger than DPHMIN-SIM. Although there are extra CPU and disk overhead, they are not as significant as we observed from this simulation.

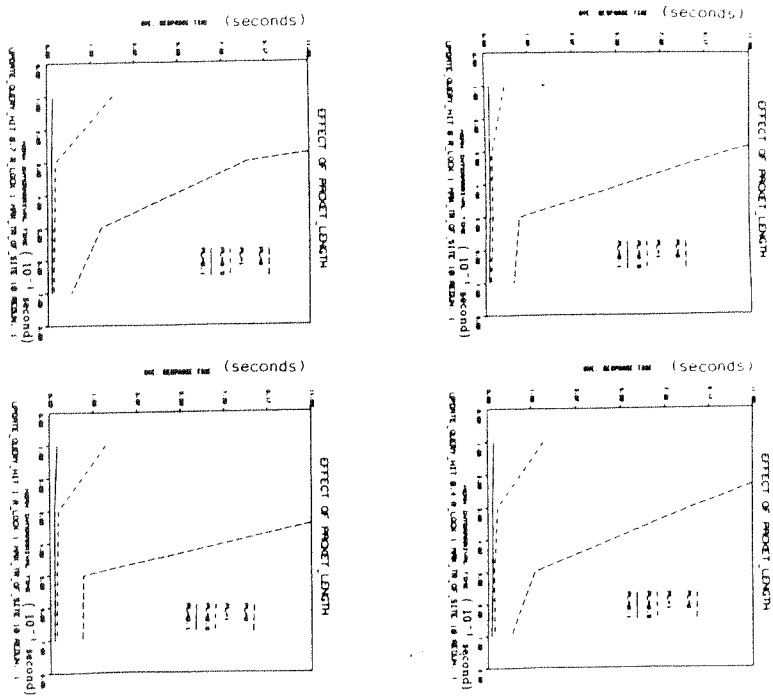


Figure 6.12 The effect of packet size to the mean response time.

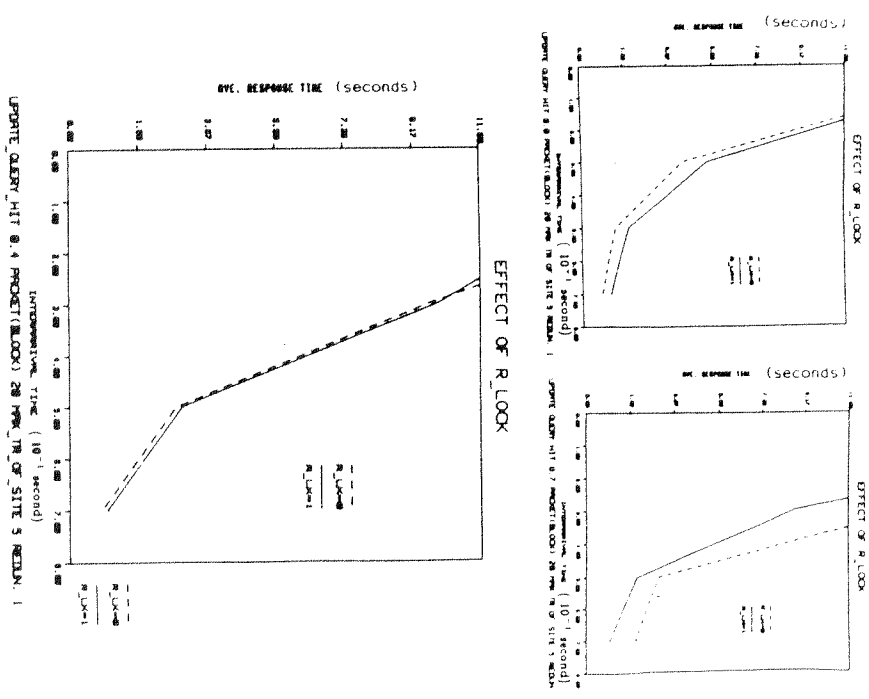


Figure 6.13 The effect of read_lock releasing time to the mean response time.

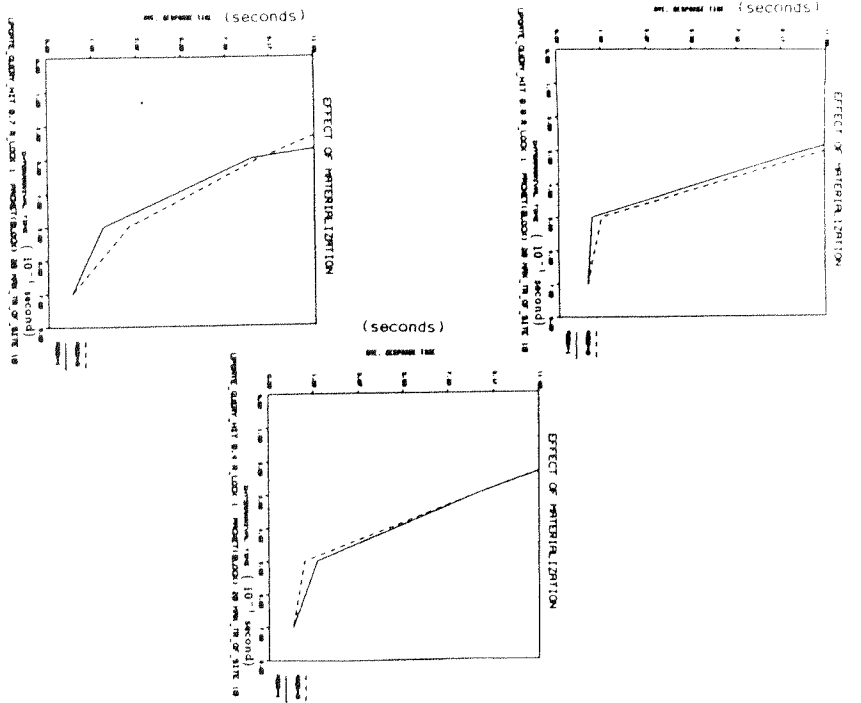


Figure 6.14 The effect of materialisation to the mean response time.

In this chapter, we have developed a queuing model to simulate DDBLMN on a UNIX system. The simulator DDBLMN-SIM implements the queuing model using DDBLMN-SMPL, which is a queue management package embedded in DDBLMN-SIM. DDBLMN was simulated under various conditions and control strategies. The statistical results show that DDBLMN is feasible and correct. The relative performance of various control strategies are compared in the simulations. The results are consistent with the expectation projected from the qualitative analysis.

The absolute values obtained from the simulation may not have significant meaning to a system designer who wants to adapt the simulation results into a real system. However, the relative performance among various control strategies would be a precious reference.

CHAPTER VII

CONCLUSIONS

7.1 Summary and Conclusions

In this research, we studied the design of a distributed database on a local computer system connected by a multiaccess network. Three key issues are addressed: data distribution, distributed query processing, and concurrency control. Based on the characteristics of local multiaccess networks, efficient strategies are proposed. The system is simulated in an event-driven queuing model to show the correctness, feasibility, and efficiency.

Due to the broadcast capability of multiaccess networks, each update can be processed at a cost independent of the allocation and the number of copies of files to be allocated. As a result, the file allocation problem, which is a subproblem of data distribution, is simplified from the problem of "where to locate copies of files?" on general networks to the problem of "should a copy of the file be located in a particular site?" Some special cases are reduced from NP-hard on general networks to polynomially solvable on local multiaccess networks. Although the general problem remains NP-hard, the complexity is reduced by decomposing the problem into N^{2m} knapsack problems. Since there exists fully polynomial approximation algorithms for the knapsack problem, efficient heuristic solutions can be easily developed based on this framework.

A transaction processing model that integrates query processing and concurrency control is proposed in the thesis. In identifying the target data in query processing, redundant materialization and semi-joins with complement transmissions are adapted. Redundant materialization helps reduce the size of intermediate results during the target data identification. On local multiaccess networks, a broadcast attribute can be semi-joined with attributes at all sites simultaneously, hence reducing the amount of data transfer as compared to that of general networks. The broadcast network allows efficient distribution of status information and permits queries to be scheduled dynamically based on the current status of the database. Although suboptimal, dynamic query processing eliminates the errors incurred in estimating status information in static algorithms, and has been found to perform better than static algorithms. An extremum-identification algorithm is proposed to aid the dissemination of status information. One of the potential directions for future research on query processing is intelligent query processing, in which knowledge is collected and extracted during the system operation and employed in later query processing. Due to the prematurity of the distributed databases, there is no sufficient knowledge available in this field to design a meaningful intelligent system. Only a model is proposed for intelligent query processing in this thesis, and detailed solutions are opened to future research.

This thesis also studies the query processing problem in future communication networks, in which the communication cost between any pair of nodes only depends on a fixed constant cost. The volume of transferred data is no longer a dominating factor to the total communication cost. Instead, the cost of communication channel usage dominates the overall cost. The DQP problem can be formulated into the generalized traveling salesman problem in

this case. A number of solution methods are considered. Among them, the most promising one that can be applied in the near future is to transform the problem into the standard traveling salesman problem and then solve it with heuristic algorithms. This strategy was evaluated. Results show that the algorithm is efficient with insignificant sacrifice to accuracy.

In contrast to conventional approach, our proposed concurrency control algorithm is embedded into dynamic query processing such that potential blocking due to unavailability of queried relations can be avoided. As a result of using multiaccess/broadcast networks, an efficient, serializable, and deadlock-free concurrency control algorithm is developed in the thesis. The broadcast capability allows locks to be known to all sites simultaneously, and hence eliminates the use of explicit lock messages. The serial transmission property of the network allows the database to be driven by network events. Therefore, the network itself serves as a synchronization mechanism.

Finally, the system is simulated based on a queuing network model. The correctness and the feasibility are demonstrated by the simulations. Furthermore, different strategies under various conditions are also compared by simulations.

In summary, the use of local multiaccess/broadcast networks not only simplifies the design of distributed databases, but also reduces the processing overhead. As a consequence, the distributed approach for a database in a local environment becomes promising as the technology of small computer systems and the local area networks advances.

7.2 Suggestions For Future Research

There remains problems to be studied in the future. These include the problem of recovery when the network or a site fails, the reliable broadcast of messages, the use of multiple broadcast buses or a bus with high bandwidth, the necessary modification on the operating system to support real-time processing of information received from the bus, the dynamic file allocation problem (file migration problem), the integration of file allocation with transaction processing, better solutions for GTSP that models future high speed communication networks, and the intelligent query processing. Some of them are discussed in the rest of this chapter.

7.2.1 High Bandwidth and Clustered Local Multiaccess Networks

The most serious drawback of a multiaccess/broadcast network with a single bus is that only serial transmission is allowed on the bus. Using current communication technologies, the bus may become the bottleneck of the system when network activities increase. One way to enhance the transmission capability is to use high speed transmission medium, such as optical fibers, with enhanced network interfaces. Another way is to use multiple bus networks. Both approaches retain the broadcast capability. Therefore, there may be only a little effect on our proposed concurrency control algorithms. For the query processing problem, the existing algorithms may have to be changed since the overhead of local processing is more significant. As discussed in Chapter IV, the semi-join approach may no longer be adequate in query processing. Totally new strategies may be needed. For the file allocation problem, the integer programming approach may remain adequate, but the solutions are much more complicated.

As the number of computer systems and other devices in the network grows, a single multiaccess network may not be sufficient to support the entire system. It is not unusual that the computing facilities owned by an organization are made by different vendors. The systems made by a particular vendor are grouped as a cluster and interconnected by a distinct vendor-supported network with its own communication protocols. These "local" networks are interconnected by one or more backbone networks as a logically integrated network. The Andrew system of the Carnegie-Mellon University is an example of such systems [MOR86]. Currently, Andrew consists of 17 Ethernets and 2 ProNet rings linked by fiber-optic cables. Over 600 computers are supported by the network currently, and over 5000 workstations are to be supported in the future. It is perceivable that this type of frameworks will become one of the most popular computing environments in the near future due to the success of the Andrew project.

Although Ethernet type networks seem to be the most attractive alternative for each individual network, the characteristics of the network as a whole may not be the same as that of a single Ethernet network. For example, the communication overhead across different networks is obviously higher than that within a local network. Different control strategies are required for a DDDBMS on this type of environment. For example, the overhead of remote file access and update is no longer independent of the locations of files. Each local network may have independent communication such that communications can proceed in all local networks simultaneously. The network no longer serves as a synchronization mechanism for concurrency control. Furthermore, new techniques to collect statistics are needed to support dynamic query processing.

7.2.2 Integrated Design of FAP and PQP

In general, the file allocation and query processing problems are interdependent of each other. To process a query, the allocation of files should be known. Likewise, to allocate files, the query access pattern and query processing strategies should also be known. Currently, only the allocation of files is considered in query processing, while query processing strategies and query patterns are not considered in the file allocation. Also, each query is assumed to be accessing a single file. Solving query processing and file allocation problems together is extremely complex, especially when delay and availability constraints are also included. Apers proposed a solution for general networks [AP82]. The problem is solved as follows. First, the query processing order is optimized independently for each query or update using distinct copies of files. Distinct copies are located at virtual sites and the file allocation problem is the mapping of the virtual sites to physical sites such that the total communication costs are minimal. This framework may not be the best approach for a DDDBMS, which uses dynamic scheduling and redundant materialization. In allocating files for a system with dynamic scheduling, it is really difficult to predict the queries a system may take during query processing. As a result, the schedules that the file allocation relied upon may be quite different from the real schedule that the query scheduler generates.

One of the key factors to reduce the communication overhead in query processing is to store as many related relations as possible at a site, where the related relations are the relations to be joined. One of the better approaches we are considering is to impose such factors into the file allocation algorithm such that it tends to allocate queried relations by a transaction into as few sites as possible, in addition to allocating a file to sites that frequently access it. We

believe that this approach will be a good direction for future research to integrate PAP and DQP together.

7.2.3 File Migration problem

The problem of file migration arises from the dynamic nature of queries. Although there may be a locality of access for a file, there are occasionally very few activities inside the locality, and the file is accessed outside its locality. Moreover, the locality of access may be time-varying, as in global networks distributed over different time zones. It is efficient to allow the file to migrate to major access sites.

A typical migration method is to examine the allocation periodically and to reallocate a file if necessary. Over a period of t intervals, there are a minimum of 2^t possible alternatives. The problem of determining the best time to execute a migration with minimum overhead is NP-hard on general network. The problem complexity on DDM-MIN may be reduced, however, efficient solutions are still opened to future research.

7.2.4 Intelligent Query Processing

As we mentioned in Section 7.1, the detailed design of intelligent query processing is still open. The problem includes the representation, acquisition and abstraction of knowledge, learning module design, expert scheduler design, knowledge base design, as well as the interactions among system components. The major difficulty is on the knowledge acquisition due to the lack of expertise in this field. Comparing to other intelligent systems, we do not have good knowledge on query processing as the internal database activities are transparent implicitly to the users. Thus, a powerful learning model is the key

design issue of such intelligent query processing systems.

7.2.5 Generalized Traveling Salesman Problem

The performance of the TSP-transformation approach for the GTSP problem is obviously bounded by the performance of the TSP problem. Furthermore, the problem size after transformation is expanded to $N=3n+m$, where n is the number of nodes and m is the number of clusters. For a large problem, even a complexity of $O(N^2)$ could be substantial. In this sense, a direct-solving solution without transformation may be a good candidate to provide better solutions. One of the potential approaches we are studying is the "decomposition" approach, in which the GTSP problem is decomposed into several sub-problems. An example algorithm in this approach is shown as follows.

Example GTSP Algorithm 1

- (a) Select an arbitrary node from each cluster.
- (b) Connect all selected nodes together as a TSP tour.
- (c) Refine the solution.

There are many ways to select nodes from clusters in step (a). Obviously, the center node, which is the node closest to the geographic center of the cluster, is a good candidate. Another good alternative is the node closest to the geographic center of the graph.

Among various potential algorithms we will consider for the refinement process, the following *rubber-band* process is one of the promising approaches.

Rubber-Band Process

Assume the TSP path $P = (0, 1, \dots, n)$ is produced in step (b) and node i is

the node selected from cluster i in step (a). Consecutive nodes on the TSP path is numbered consecutively, and node 0 is assumed to be the node closest to the geographic center.

```

For all  $k$  in  $(0, 1, \dots, n)$  do
  set  $k_1 = (k+1) \bmod (n+1)$ 
  set  $k_2 = (k+2) \bmod (n+1)$ 
  draw a straight line between nodes  $k$  and  $k_2$ 
  from cluster  $k+1$ , select node  $k_1'$  which is closest to the line
  replace node  $k_1$  in  $P$  with  $k_1'$ 
done

```

The solution can be further refined as many times as needed. Such a process is analogous to the following process:

"Consider each node as a pole, a rubber band is used to circle at least one pole from each cluster with the minimum length."

It is obvious that a necessary condition for an optimal solution is that the rubber band should touch the pole closest to the line drawn between the two poles selected from two neighboring clusters.

Since the algorithms are heuristic, substantial experimentation is required to study the performance of different algorithms under different conditions.

Currently, there is no way to draw any conclusion from the above discussion. However, this could be a good direction for future research.

LIST OF REFERENCES

LIST OF REFERENCES

- [AHO74] A. V. Aho, J. E. Hopcroft, and J. D. Ullman, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Mass., 1974, pp. 199-200.
- [APE82] P. M. G. Apers, "Centralized or Decentralized Data Allocation," *Distributed Data Sharing Systems*, R. P. Van de Riet and W. Litwin, eds., North-Holland, New York, 1982.
- [APE83] P. M. G. Apers, A. R. Heyner, and S. B. Yao, "Optimization Algorithms for Distributed Queries," *IEEE Trans. on Software Engineering*, Vol. SE-9, No. 1, Jan. 1983, pp. 57-68.
- [BAN79] J.-S. Banino, C. Kaiser, and H. Zimmermann, "Synchronization for Distributed Systems using a Single Broadcast Channel," *Proceedings of the 1st Int'l Conf. On Distributed Computing Systems*, Huntsville, Oct. 1979, pp. 330-338.
- [BEA59] J. Bearwood, J. H. Halton and J. M. Hammersley, "The Shortest Route Through Many Points," *Proceedings of the Cambridge Philosophical Society* 55, 1959, pp. 299-327.
- [BEL62] R. Bellman, "Dynamic Programming Treatment of the Traveling Salesman Problem," *JACM*, Vol. 9, No. 1, 1962, pp. 61-63.
- [BEL68] M. Bellmore and G. L. Nemhauser, "The Traveling Salesman Problem: A Survey," *Operations Research*, Vol. 16, 1968, pp. 538-558.
- [BEL71] M. Bellmore and J. C. Malone, "Pathology of Traveling Salesman Subtour-Elimination Algorithms," *Operations Research*, Vol. 19, 1971, pp. 278-307.
- [BEN83] J. L. Bentley, "A Case Study in Applied Algorithm Design: The Traveling Salesman Problem," *Dept. of Comp. Sci. and Math., Carnegie-Mellon Univ.*, April 7, 1983.
- [BER81a] P. A. Bernstein and D. M. Chin, "Using Semi-Joins to Solve Relational Queries," *JACM*, Vol. 28, No. 1, Jan. 1981, pp. 25-40.
- [BER81b] P. A. Bernstein and N. Goodman, "Concurrency Control in Distributed Database Systems," *ACM Computing Surveys*, Vol. 13, No. 2, June 1981, pp. 185-221.
- [BER81c] P. A. Bernstein, N. Goodman, E. Wong, C. L. Reeve, and J. B. Rothnie, "Query Processing in a System for Distributed Databases (SDD-1)," *ACM Trans. on Database Systems*, Vol. 6, No. 4, Dec. 1981, pp. 602-625.
- [BHA82a] B. Bhargava and C. T. Hua, "Classes of Serializable Histories and Synchronization Algorithms in Distributed Database Systems," *Proceedings of IEEE International Conference on Distributed Computing Systems*, Miami, Florida, Oct. 1982, pp. 438-446.
- [BHA82b] B. Bhargava, "Performance Evaluation of the Optimistic Approach to Distributed Database Systems and its Comparison," *Proceedings of IEEE International Conference on Distributed Computing Systems*, Miami, Florida, Oct. 1982, pp. 508-517.
- [BHA83] B. Bhargava and C. T. Hua, "A causal Model for Analyzing Distributed Concurrency Control Algorithms," *IEEE Trans. On Software Engineering*, Vol. SE-9, NO. 4, July 1983, pp 470-486.

- [BHR82] A. Birrell, "Grapevine: An Exercise In Distributed Computing," *CACM*, Vol. 15, No. 4, April 1982, pp. 260-273.
- [BUR79] R. E. Burkard, "Traveling Salesman and Assignment Problems: A Survey," *Annals of Discrete Mathematics*, Vol. 4, pp. 193-215, North-Holland, 1979
- [CER82] S. Ceri, G. Paolini, G. Pelagatti and F. A. Schreiber, "Distributed Database Research at the Politecnico Milano," *IEEE Data Engineering*, Vol. 5, No. 4, Dec. 1982, pp. 9-13.
- [CER83] S. Ceri, S. Navathe, and G. Wiederhold, "Distribution Design of Logical Database Schemas," *IEEE Trans. on Software Engineering*, Vol. SE-9, No. 4, July 1983, pp. 487-504.
- [CER84] Stefano Ceri and Giuseppe Pelagatti, "Distributed Databases: Principles and Systems," McGraw-Hill, New York, 1984.
- [CHA83] J. M. Chang, "LAMBDA: A Distributed Database System," Bell Telephone Lab. Tech. Report, March 1983.
- [CHA84a] J. S. Chang, "Simplifying Distributed Database Systems Design by Using a Broadcast Network," *Proc. ACM SIGMOD Int'l Conf. on Management of Data*, June 1984, pp. 223-233.
- [CHA84b] J. M. Chang and N. F. Maxemchuk, "Reliable Broadcast Protocols," *ACM Trans. on Computer Systems*, Vol. 2, No. 3, August 1984, pp. 251-273.
- [CHA85] J. S. Chang, "LAMBDA: A Distributed Database System for Local Area Networks," *Database Engineering*, Vol. 8, No. 2, June 1985, pp. 76-83.

- [CHR84] A. L. P. Chen, and V. O. K. Li, "Improvement Algorithms for Semijoin Query Processing Programs in Distributed Database Systems," *IEEE Trans. on Computers*, Vol. C-33, No. 11, Nov. 1984, pp. 959-967.
- [CHI84] D. M. Chiu, P. A. Bernstein, and Y. C. Ho, "Optimizing Chain Queries in a Distributed Database System," *SIAM J. of Comp.*, Vol. 13, No. 1, Feb. 1984, pp. 116-134.
- [CHR72] N. Christofides and S. Eilon, "Algorithms for Large-scale Travelling Salesman Problems," *Operational Research Quarterly*, Vol. 23, No. 4, 1972, pp. 511-518.
- [CHU69] W. W. Chu, "Multiple File Allocation in a Multiple Computer System," *IEEE Trans. on Computers*, Vol. C-18, No. 10, Oct. 1969, pp. 885-889.
- [CHU82] W. W. Chu, and P. Hurley, "Optimal Query Processing For Distributed Database Systems," *IEEE Trans. on Computers*, Vol. C-31, No. 9, Sept. 1982, pp. 835-850.
- [CLA78] A. Claus, "A Simultaneous Enumeration Approach to the Traveling Salesman Problem," *Studies In Applied Math.*, Vol. 58, 1978, pp. 159-163.
- [COD70] E. F. Codd, "A Relation Model of Data for Large Shared Data Banks," *CACM*, Vol. 13, 1970, pp. 377-387.
- [DAN54] G. B. Dantzig, D. R. Fulkerson, and S. Johnson, "Solution of a Large-Scale Traveling-Salesman Problem," *Operations Research*, Vol. 2, 1954, pp. 393-410.

- [DOW82] L. W. Dowdy and D. V. Foster, "Comparative Models of the File Assignment Problem," *ACM Computing Surveys*, Vol. 14, No. 2, June 1982, pp. 287-313.
- [ESW74] K. P. Eswaran, "Placement of Records in a File and File Allocation in a Computer Network," *Information Processing 74*, IFIPS, North Holland Publishing Co., 1974.
- [FIS78] G. S. Fishman, "Principles of Discrete Event Simulation," *Wiley-Interscience, 1978. Management Science*, Vol. 18, No. 9, May 1972, pp. 465-491.
- [GEO72] A. M. Geoffion and R. E. Marsten, "Integer Programming: A Framework and State-of-art Survey," *Management Sci.*, Vol. 18, No. 9, pp. 465-491, May 1972.
- [GIL64] P. C. Gilmore and R. E. Gomory, "Sequencing A One State-Variable Machine: A Solvable Case of the Traveling Salesman Problem," *Operations Research*, Vol. 12, Sep.-Oct., 1964, pp. 655-679.
- [GON62] *Solution to the Traveling Salesman Problem by Dynamic Programming on the Hypercube*, Tech. Rep. No. 18, O.R. Center, M.I.T., 1962.
- [GOU81] M. G. Gouda and U. Dayal, "Optimal Semi-join Schedules for Query Processing in Local Distributed Database Systems," *Proc. ACM SIGMOD Conference*, May 1981, pp. 164-175.
- [GUP68] J. N. D. Gupta, "Traveling Salesman Problem-A Survey of Theoretical Developments and Applications," *Operach*, Vol. 5, 1968, pp. 181-192.

- [GUP78] J. N. D. Gupta, "A Search Algorithm for the Traveling Salesman Problem," *Compu. & Ops. Res.*, Vol. 5, 1978, pp. 243-250.
- [HAR62] W. W. Hardgrave and G. L. Nemhauser, "On the Relation Between the Traveling-Salesman and the Longest-Path Problems," *Operations Research*, Vol. 10, No. 5, Jan.-Feb., 1962, pp. 647-657.
- [HEL62] M. Held, and R. M. Karp, "A Dynamic Programming Approach to Sequencing Problems," *SIAM*, Vol. 10, No. 1, 1962, pp. 196-210.
- [HEL70] M. Held, and R. M. Karp, "The Traveling Salesman Problem and Minimum Spanning Trees," *Operations Research*, Vol. 18, No. 6, 1970, pp. 1138-1162.
- [HEL71] M. Held and R. M. Karp, "The Traveling Salesman Problem and Minimum Spanning Trees: Part II," *Math. Prog.*, Vol. 1, pp. 6-25, North-Holland, 1971.
- [HEV79a] A. R. Heyner, *Query Processing in Distributed Database Systems*, Ph.D. Dissertation, Computer Sciences Department, Purdue University, CS79-4, Dec. 1979.
- [HEV79b] A. R. Heyner, "The Optimization Of Query Processing On Distributed Database Systems," *IEEE Trans. on Software Engineering*, Vol. SE-5, No. 3, May 1979, pp. 177-187.
- [HEV85] Alan R. Heyner, O. Qi Wu, and S. B. Yao, "Query Optimization on Local Area Networks," *ACM Transaction on Office Information*, Vol. 3, No. 1, Jan. 1985, pp. 35-62.
- [HON76] S. Hong, "The Asymmetric Traveling Salesman Problem and Spanning Arborescences of Directed Graphs," *Proceedings of Symposium on Operations Research*, Heidelberg, Nov. 1976.

- [HOU80] D. J. Houck Jr., J. C. Byard, M. Queyranne, and R. R. Vennganti, "The Traveling Salesman Problem As A Constrained Shortest Path Problem: Theory and Computational Experience," *Operach*, Vol. 17, No. 2&3, 1980, pp. 93-109.
- [IEEE83] Institute for Electrical and Electronic Engineers, "IEEE project 802, Local Area Network Standards," 1983.
- [JUWA84] J. Y. Juang and B. W. Wah, "Unified Window Protocols for Contention Resolution in Local Multiaccess Networks," *Proc. of 1984 IEEE INFOCOM*, San Francisco, CA, April 1984, PP. 97-104.
- [KAR64] R. L. Karg and G. L. Thompson, "A Heuristic Approach to Solving Travelling Salesman Problems," *Management Science*, Vol. 10, No. 2, Jan., 1964, pp. 225-248.
- [KAR79] R. M. Karp, "A Patching Algorithm for the Nonsymmetric Traveling Salesman Problem," *SIAM J. Computing*, Vol. 8, No. 4, Nov., 1979, pp. 561-573.
- [KER82] Larry Kerschberg, Peter D. Ting, and S. Bing Yao, "Query Optimization in Star Computer Networks," *ACM Transactions On Database Systems*, Vol. 7, No. 4, Dec. 1982, pp. 678-711.
- [KRIS84] R. Krishnamurthy and S. P. Morgan, "Distributed Query Optimization: An Engineering Approach," *Proc. Int'l Conf. on Data Engineering*, Los Angeles, CA, April 1984, pp. 220-227.
- [KUM70] S. Kumar, "Optimal Path through k Specified Sets of Nodes," *Indian Journal of Mathematics*, Vol. 12, No. 1, pp. 25-30, Jan. 1970.
- [KUN81] H. T. Kung and J. T. Robinson, "Optimistic Methods for Concurrency Control," *ACM Trans. On Database Systems*, Vol. 6, No. 2, pp. 213-226, Jan. 1981.

- [LIN65] S. Lin, "Computer Solutions of the Traveling Salesman Problem," *Bell System Tech. Journal*, Vol. 44, 1965, pp. 2245-2269.
- [LIN84] B. G. Lindsey, L. M. Haas, C. Mohan, P. F. Wilms, and R. A. Yost, "Computation and Communication in R," *ACM Trans. On Computer Systems*, Vol. 2, No. 1, Feb. 1984, pp. 24-38.
- [LIT63] J. D. C. Little, K. G. Murty, D. W. Sweeney, and G. Karel, "An Algorithm for the Traveling Salesman Problem," *Operations Research*, Vol. 11, Nov.-Dec., 1963, pp. 972-989.
- [LU85] Hongjun Lu, *Distributed Query Processing With Load Balancing In Local Area networks*, Ph.D. Dissertation, Computer Sciences Department, University of Wisconsin-Madison, Computer Sciences technical Report #624, 1985.
- [MAC80] M. H. MacDougall, "SMPL - A Simple Portable Simulation Language," *Amdahl Technical Report*, April, 1980.
- [MET76] R. F. Metcalfe and D. R. Boggs, "Ethernet: Distributed packet switching for local computer networks," *CACM*, Vol. 19, No. 7, July 1976, pp. 395-404.
- [MIL76] P. Milhous, "Integer Programming Approaches to the Travelling Salesman Problem," *Math. Prog.*, Vol. 10, 1976, pp. 367-378.
- [MIL78] P. Milhous, "Using Cutting Planes to Solve the Symmetric Traveling Salesman Problem," *Math. Prog.*, Vol. 15, pp. 177-188, North-Holland, 1978.
- [MOR86] J. H. Morris, M. Satyanarayanan, M. H. Conner, J. H. Howard, D. S. H. Rosenthal, and F. D. Smith, "Andrew: A Distributed Personal Computing Environment," *CACM*, Vol. 29, No. 3, Mar. 1986, pp. 184-201.

- [NGU181a] Gia Toan Nguyen and G. Sergeant, "Distributed Architecture and Decentralized Control for a Local Network Database System," *Proceedings of ACM International Computing Symposium*, London, March 1981.
- [NGU81b] Gia Toan Nguyen, "Distributed Query Management for a Local Network Database System," *Proceedings of the Second International Conference on Distributed Computing Systems*, Paris, April 1981.
- [PAN64] S. N. N. Pandit, "Some Observations On the Longest Path Problem," *Operations Research*, Vol. 12, No. 2, 1964, pp. 361-364.
- [PAP79] C. H. Papadimitriou, "The Serializability of Concurrent Database Updates," *JACM*, Vol. 26, No. 4, Oct. 1979, pp. 631-653.
- [PAP82] C. H. Papadimitriou, "A Theorem in Database Concurrency Control," *JACM*, Vol. 29, No. 4, Oct. 1982, pp. 998-1006.
- [PEE78] Richard Peebles and Eric Manning, "System Architecture for Distributed Data Management," *IEEE Computer*, Vol. 11, No. 1, Jan. 1978, pp. 40-47.
- [POP83] G. J. Popok and G. Thiel, "Distributed Data Management Issues in the LOCUS System," *IEEE Data Engineering*, Vol. 6, No. 2, June 1983, pp. 63-68.
- [RAM83] C. V. Ramamoorthy and B. W. Wah, "The Isomorphism of Simple File Allocation," *IEEE Trans. On Computers*, Vol. C-32, No. 3, March 1983, pp. 221-232.
- [ROT80] J. Rohnie, P. A. Bernstein, S. Fox, N. Goodman, M. Hammer, T. A. Landers, C. Reeve, D. W. Shipman, and E. Wong, "Introduction to a System for Distributed Databases (SDD-1)," *ACM Trans. On Database Systems*, Vol. 5, No. 1, March 1980, pp. 1-17.
- [SAC84] G. M. Sacco, "Distributed Query Evaluation in Local Area Networks," *Proc. of Int'l Conf. On Data Engineering*, Los Angeles, CA, April 1984, pp. 510-516.
- [SHO82] J. H. Shoch, Y. K. Datal, D. D. Redell, and R. C. Crane, "Evolution of the Ethernet Local Computer Network," *IEEE Computer*, Vol. 15, No. 8, Aug. 1982, pp. 10-27.
- [SRI70] S. S. Srivastava, S. Kumar, R. C. Garg and P. Sen, "Generalized Traveling-Salesman Problem through n Sets of Nodes," *Journal of the Canadian Operational Research Society*, pp. 97-101, 1970.
- [STA84] William Stallings, "Local Networks," *ACM Computing Surveys*, Vol. 16, No. 1, Mar. 1984, pp. 3-42.
- [STA85] William Stallings, *Data And Computer Communications*, Macmillan, 1985.
- [STO77] M. Stonebraker and E. Neuhold, "A distributed Data Base Version of INGRES," *Proc. 1977 Berkeley Workshop on Distributed Data Management and Computer Networks*, May 1977, pp. 19-36.
- [STU80] H. Sturgis, J. Michell, and J. Israel, "Issues In The Design And Use Of A Distributed File System," *ACM Operating System Review*, Vol. 14, No. 3, July 1980, pp. 55-69.
- [SYK80] David J. Sykes, "The Economics Of Distributed Systems," *Digest of Papers, COMPCOM*, Fall 1980, pp. 8-15.
- [THO79] R. H. Thomas, "A Majority Consensus Approach for Concurrency Control," *ACM Trans. On Database Systems*, Vol. 4, No. 2, June 1979, pp. 180-209.

- [ULL82] J. D. Ullman, *Principles of Database Systems*, Computer Science Press, 1982.
- [WAH83] B. W. Wah and J. Y. Jiang, "An Efficient Protocol for Load Balancing on CSMA/CD Networks," *Proc. 8th Conference on Local Computer Networks*, Minneapolis, MN, Oct. 1983, pp. 55-61.
- [WAH84a] B. W. Wah, "File Placement on Distributed Computer Systems," *IEEE Computer*, Vol. 17, No. 1, Jan. 1984, pp. 23-32.
- [WAH84b] B. W. Wah and Y. N. Lien, "The File-Assignment And Query-Processing Problems In Local Multiaccess Networks," *Proc. of Intl Conf. On Data Engineering*, Los Angeles, CA, April 1984, pp. 228-235.
- [WAH85] B. W. Wah and Y. N. Lien, "Design of Distributed Databases On Local Computer Systems With A Multiaccess Network," *IEEE Trans. On Software Eng.*, Vol. SE-11, No. 7, July 1985, pp. 606-619.
- [WAH85b] B. W. Wah and J. Y. Juang, "Resource Scheduling for Local Computer Systems with a Multiaccess Network," *IEEE Trans. on Computers*, Vol. C-34, No. 12, Dec. 1985, pp. 1144-1157.
- [WON77] E. Wong, "Retrieving Dispersed Data From SDD-1: A System for Distributed Databases," *Proc. 2nd Berkeley Workshop on Distributed Data Management and Computer Networks*, May, 1977, pp. 217-235.
- [YAO77] S. B. Yao, "Approximating Block Accesses In Database Organization," *CACM*, Vol. 20, No. 4, April 1977, pp. 260-261.
- [YAO79] S. B. Yao, "Optimization of Query Evaluation Algorithms," *ACM Trans. On Database Systems*, Vol. 4, No. 2, June 1979, pp. 133-155.

- [YU82] C. T. Yu, C. C. Chang, and Y. C. Chang, "Two Surprising Results in Processing Simple Queries in Distributed Databases," *Proc. COMPSAC 82*, Nov. 82, pp. 377-384.
- [YU83] C. T. Yu and C. C. Chang, "On The Design of A Query Processing Strategy In A Distributed Database Environment," *Proc. SIGMOD Conference*, Vol. 13, No. 4, 1983, pp. 30-39.
- [YU84a] C. T. Yu and C. C. Chang, "Distributed Query Processing," *ACM Computing Surveys*, Vol. 16, No. 4, Dec. 1984, pp. 399-433.
- [YU84b] C. T. Yu, Z. M. Ozsoyoglu, and C. C. Chang, "Optimization of Distributed Tree Queries," *Journal of Computer and System Sciences*, Vol. 29, No. 3, Dec. 1984, pp. 409-445.

APPENDIX

```
1 /******
2 NAME ..... DOBLAN-SIW simulator
3 PURPOSE ..... To simulate a DBMS on a local multiaccess network
4
5 REFERENCE ..... B. B. Sridhar, N. Lind,
6 ..... Design and Distribution of Database On Local Computer
7 ..... Systems With A Multiprocessor Architecture,
8 ..... IEEE Trans. On Software Engineering,
9 ..... Vol. SE-11, No. 7, July 1985, pp. 600-619
10
11 STD INPUT ..... seed seed to random number generator ( int )
12 ..... length of execution time ( second )
13 ..... mean transaction arrival time ( 0.1 second )
14 ..... transaction of updated transaction (librv: 0-10)
15 ..... time to release read lock ( 0 SECONDS )
16
17 STD OUTPUT ..... DOBLAN SMPLE information
18 ..... online containing following information:
19 ..... 1. simulation execution time at the end of
20 ..... a simulation execution
21 ..... 2. average transaction processing time in second
22 ..... 3. average (ratio of the time spent on bus to the total
23 ..... average (ratio of time of TR ) on total transactions
24 ..... processing time of TR ) on total transactions
25 ..... 4. average (ratio of the time spent on CPU to the total
26 ..... lost processing time of TR ) on total transactions
27 ..... 5. utilization of disk at site 1
28 ..... 6. average queue length of disk at site 1
29 ..... 7. utilization of disk at site 2
30 ..... 8. average queue length of disk at site 2
31 ..... 9. utilization of disk at site 3
32 ..... 10. average queue length of disk at site 3
33 ..... 11. utilization of disk at site 4
34 ..... 12. average queue length of disk at site 4
35 ..... 13. utilization of disk at site 5
36 ..... 14. average queue length of disk at site 5
37 ..... 15. utilization of disk at site 6
38 ..... 16. average queue length of disk at site 6
39 ..... 17. utilization of network interface at site 1
40 ..... 18. average queue length of network interface at site 1
41 ..... 19. utilization of network interface at site 2
42 ..... 20. average queue length of network interface at site 2
43 ..... 21. utilization of network interface at site 3
44 ..... 22. average queue length of network interface at site 3
45 ..... 23. utilization of network interface at site 4
46 ..... 24. average queue length of network interface at site 4
47 ..... 25. utilization of network interface at site 5
48 ..... 26. average queue length of network interface at site 5
49 ..... 27. utilization of network interface at site 6
50 ..... 28. average queue length of network interface at site 6
51 ..... 29. CPU time in second used in the simulation
52 ..... 30.
53 ..... 31.
54 ..... 32.
55 ..... 33.
56 ..... 34.
57 ..... 35.
58 ..... 36.
59 ..... 37.
60 ..... 38.
61 ..... 39.
62 ..... 40.
63 ..... 41.
64 ..... 42.
65 ..... 43.
66 ..... 44.
67 ..... 45.
68 ..... 46.
69 ..... 47.
70 ..... 48.
71 ..... 49.
72 ..... 50.
73 ..... 51.
74 ..... 52.
75 ..... 53.
76 ..... 54.
77 ..... 55.
78 ..... 56.
79 ..... 57.
80 ..... 58.
81 ..... 59.
82 ..... 60.
83 ..... 61.
84 ..... 62.
85 ..... 63.
86 ..... 64.
87 ..... 65.
88 ..... 66.
89 ..... 67.
90 ..... 68.
91 ..... 69.
92 ..... 70.
93 ..... 71.
94 ..... 72.
95 ..... 73.
96 ..... 74.
97 ..... 75.
98 ..... 76.
99 ..... 77.
100 ..... 78.
101 ..... 79.
102 ..... 80.
103 ..... 81.
104 ..... 82.
105 ..... 83.
106 ..... 84.
107 ..... 85.
108 ..... 86.
109 ..... 87.
110 ..... 88.
111 ..... 89.
112 ..... 90.
113 ..... 91.
114 ..... 92.
115 ..... 93.
116 ..... 94.
117 ..... 95.
118 ..... 96.
119 ..... 97.
120 ..... 98.
121 ..... 99.
122 ..... 100.
123 ..... 101.
124 ..... 102.
125 ..... 103.
126 ..... 104.
127 ..... 105.
128 ..... 106.
129 ..... 107.
130 ..... 108.
131 ..... 109.
132 ..... 110.
133 ..... 111.
134 ..... 112.
135 ..... 113.
136 ..... 114.
137 ..... 115.
138 ..... 116.
139 ..... 117.
140 ..... 118.
141 ..... 119.
142 ..... 120.
143 ..... 121.
144 ..... 122.
145 ..... 123.
146 ..... 124.
147 ..... 125.
148 ..... 126.
149 ..... 127.
150 ..... 128.
151 ..... 129.
152 ..... 130.
153 ..... 131.
154 ..... 132.
155 ..... 133.
156 ..... 134.
157 ..... 135.
158 ..... 136.
159 ..... 137.
160 ..... 138.
161 ..... 139.
162 ..... 140.
163 ..... 141.
164 ..... 142.
165 ..... 143.
166 ..... 144.
167 ..... 145.
168 ..... 146.
169 ..... 147.
170 ..... 148.
171 ..... 149.
172 ..... 150.
173 ..... 151.
174 ..... 152.
175 ..... 153.
176 ..... 154.
177 ..... 155.
178 ..... 156.
179 ..... 157.
180 ..... 158.
181 ..... 159.
182 ..... 160.
183 ..... 161.
184 ..... 162.
185 ..... 163.
186 ..... 164.
187 ..... 165.
188 ..... 166.
189 ..... 167.
190 ..... 168.
191 ..... 169.
192 ..... 170.
193 ..... 171.
194 ..... 172.
195 ..... 173.
196 ..... 174.
197 ..... 175.
198 ..... 176.
199 ..... 177.
200 ..... 178.
201 ..... 179.
202 ..... 180.
203 ..... 181.
204 ..... 182.
205 ..... 183.
206 ..... 184.
207 ..... 185.
208 ..... 186.
209 ..... 187.
210 ..... 188.
211 ..... 189.
212 ..... 190.
213 ..... 191.
214 ..... 192.
215 ..... 193.
216 ..... 194.
217 ..... 195.
218 ..... 196.
219 ..... 197.
220 ..... 198.
221 ..... 199.
222 ..... 200.
223 ..... 201.
224 ..... 202.
225 ..... 203.
226 ..... 204.
227 ..... 205.
228 ..... 206.
229 ..... 207.
230 ..... 208.
231 ..... 209.
232 ..... 210.
233 ..... 211.
234 ..... 212.
235 ..... 213.
236 ..... 214.
237 ..... 215.
238 ..... 216.
239 ..... 217.
240 ..... 218.
241 ..... 219.
242 ..... 220.
243 ..... 221.
244 ..... 222.
245 ..... 223.
246 ..... 224.
247 ..... 225.
248 ..... 226.
249 ..... 227.
250 ..... 228.
251 ..... 229.
252 ..... 230.
253 ..... 231.
254 ..... 232.
255 ..... 233.
256 ..... 234.
257 ..... 235.
258 ..... 236.
259 ..... 237.
260 ..... 238.
261 ..... 239.
262 ..... 240.
263 ..... 241.
264 ..... 242.
265 ..... 243.
266 ..... 244.
267 ..... 245.
268 ..... 246.
269 ..... 247.
270 ..... 248.
271 ..... 249.
272 ..... 250.
273 ..... 251.
274 ..... 252.
275 ..... 253.
276 ..... 254.
277 ..... 255.
278 ..... 256.
279 ..... 257.
280 ..... 258.
281 ..... 259.
282 ..... 260.
283 ..... 261.
284 ..... 262.
285 ..... 263.
286 ..... 264.
287 ..... 265.
288 ..... 266.
289 ..... 267.
290 ..... 268.
291 ..... 269.
292 ..... 270.
293 ..... 271.
294 ..... 272.
295 ..... 273.
296 ..... 274.
297 ..... 275.
298 ..... 276.
299 ..... 277.
300 ..... 278.
301 ..... 279.
302 ..... 280.
303 ..... 281.
304 ..... 282.
305 ..... 283.
306 ..... 284.
307 ..... 285.
308 ..... 286.
309 ..... 287.
310 ..... 288.
311 ..... 289.
312 ..... 290.
313 ..... 291.
314 ..... 292.
315 ..... 293.
316 ..... 294.
317 ..... 295.
318 ..... 296.
319 ..... 297.
320 ..... 298.
321 ..... 299.
322 ..... 300.
323 ..... 301.
324 ..... 302.
325 ..... 303.
326 ..... 304.
327 ..... 305.
328 ..... 306.
329 ..... 307.
330 ..... 308.
331 ..... 309.
332 ..... 310.
333 ..... 311.
334 ..... 312.
335 ..... 313.
336 ..... 314.
337 ..... 315.
338 ..... 316.
339 ..... 317.
340 ..... 318.
341 ..... 319.
342 ..... 320.
343 ..... 321.
344 ..... 322.
345 ..... 323.
346 ..... 324.
347 ..... 325.
348 ..... 326.
349 ..... 327.
350 ..... 328.
351 ..... 329.
352 ..... 330.
353 ..... 331.
354 ..... 332.
355 ..... 333.
356 ..... 334.
357 ..... 335.
358 ..... 336.
359 ..... 337.
360 ..... 338.
361 ..... 339.
362 ..... 340.
363 ..... 341.
364 ..... 342.
365 ..... 343.
366 ..... 344.
367 ..... 345.
368 ..... 346.
369 ..... 347.
370 ..... 348.
371 ..... 349.
372 ..... 350.
373 ..... 351.
374 ..... 352.
375 ..... 353.
376 ..... 354.
377 ..... 355.
378 ..... 356.
379 ..... 357.
380 ..... 358.
381 ..... 359.
382 ..... 360.
383 ..... 361.
384 ..... 362.
385 ..... 363.
386 ..... 364.
387 ..... 365.
388 ..... 366.
389 ..... 367.
390 ..... 368.
391 ..... 369.
392 ..... 370.
393 ..... 371.
394 ..... 372.
395 ..... 373.
396 ..... 374.
397 ..... 375.
398 ..... 376.
399 ..... 377.
400 ..... 378.
401 ..... 379.
402 ..... 380.
403 ..... 381.
404 ..... 382.
405 ..... 383.
406 ..... 384.
407 ..... 385.
408 ..... 386.
409 ..... 387.
410 ..... 388.
411 ..... 389.
412 ..... 390.
413 ..... 391.
414 ..... 392.
415 ..... 393.
416 ..... 394.
417 ..... 395.
418 ..... 396.
419 ..... 397.
420 ..... 398.
421 ..... 399.
422 ..... 400.
423 ..... 401.
424 ..... 402.
425 ..... 403.
426 ..... 404.
427 ..... 405.
428 ..... 406.
429 ..... 407.
430 ..... 408.
431 ..... 409.
432 ..... 410.
433 ..... 411.
434 ..... 412.
435 ..... 413.
436 ..... 414.
437 ..... 415.
438 ..... 416.
439 ..... 417.
440 ..... 418.
441 ..... 419.
442 ..... 420.
443 ..... 421.
444 ..... 422.
445 ..... 423.
446 ..... 424.
447 ..... 425.
448 ..... 426.
449 ..... 427.
450 ..... 428.
451 ..... 429.
452 ..... 430.
453 ..... 431.
454 ..... 432.
455 ..... 433.
456 ..... 434.
457 ..... 435.
458 ..... 436.
459 ..... 437.
460 ..... 438.
461 ..... 439.
462 ..... 440.
463 ..... 441.
464 ..... 442.
465 ..... 443.
466 ..... 444.
467 ..... 445.
468 ..... 446.
469 ..... 447.
470 ..... 448.
471 ..... 449.
472 ..... 450.
473 ..... 451.
474 ..... 452.
475 ..... 453.
476 ..... 454.
477 ..... 455.
478 ..... 456.
479 ..... 457.
480 ..... 458.
481 ..... 459.
482 ..... 460.
483 ..... 461.
484 ..... 462.
485 ..... 463.
486 ..... 464.
487 ..... 465.
488 ..... 466.
489 ..... 467.
490 ..... 468.
491 ..... 469.
492 ..... 470.
493 ..... 471.
494 ..... 472.
495 ..... 473.
496 ..... 474.
497 ..... 475.
498 ..... 476.
499 ..... 477.
500 ..... 478.
501 ..... 479.
502 ..... 480.
503 ..... 481.
504 ..... 482.
505 ..... 483.
506 ..... 484.
507 ..... 485.
508 ..... 486.
509 ..... 487.
510 ..... 488.
511 ..... 489.
512 ..... 490.
513 ..... 491.
514 ..... 492.
515 ..... 493.
516 ..... 494.
517 ..... 495.
518 ..... 496.
519 ..... 497.
520 ..... 498.
521 ..... 499.
522 ..... 500.
523 ..... 501.
524 ..... 502.
525 ..... 503.
526 ..... 504.
527 ..... 505.
528 ..... 506.
529 ..... 507.
530 ..... 508.
531 ..... 509.
532 ..... 510.
533 ..... 511.
534 ..... 512.
535 ..... 513.
536 ..... 514.
537 ..... 515.
538 ..... 516.
539 ..... 517.
540 ..... 518.
541 ..... 519.
542 ..... 520.
543 ..... 521.
544 ..... 522.
545 ..... 523.
546 ..... 524.
547 ..... 525.
548 ..... 526.
549 ..... 527.
550 ..... 528.
551 ..... 529.
552 ..... 530.
553 ..... 531.
554 ..... 532.
555 ..... 533.
556 ..... 534.
557 ..... 535.
558 ..... 536.
559 ..... 537.
560 ..... 538.
561 ..... 539.
562 ..... 540.
563 ..... 541.
564 ..... 542.
565 ..... 543.
566 ..... 544.
567 ..... 545.
568 ..... 546.
569 ..... 547.
570 ..... 548.
571 ..... 549.
572 ..... 550.
573 ..... 551.
574 ..... 552.
575 ..... 553.
576 ..... 554.
577 ..... 555.
578 ..... 556.
579 ..... 557.
580 ..... 558.
581 ..... 559.
582 ..... 560.
583 ..... 561.
584 ..... 562.
585 ..... 563.
586 ..... 564.
587 ..... 565.
588 ..... 566.
589 ..... 567.
590 ..... 568.
591 ..... 569.
592 ..... 570.
593 ..... 571.
594 ..... 572.
595 ..... 573.
596 ..... 574.
597 ..... 575.
598 ..... 576.
599 ..... 577.
600 ..... 578.
601 ..... 579.
602 ..... 580.
603 ..... 581.
604 ..... 582.
605 ..... 583.
606 ..... 584.
607 ..... 585.
608 ..... 586.
609 ..... 587.
610 ..... 588.
611 ..... 589.
612 ..... 590.
613 ..... 591.
614 ..... 592.
615 ..... 593.
616 ..... 594.
617 ..... 595.
618 ..... 596.
619 ..... 597.
620 ..... 598.
621 ..... 599.
622 ..... 600.
623 ..... 601.
624 ..... 602.
625 ..... 603.
626 ..... 604.
627 ..... 605.
628 ..... 606.
629 ..... 607.
630 ..... 608.
631 ..... 609.
632 ..... 610.
633 ..... 611.
634 ..... 612.
635 ..... 613.
636 ..... 614.
637 ..... 615.
638 ..... 616.
639 ..... 617.
640 ..... 618.
641 ..... 619.
642 ..... 620.
643 ..... 621.
644 ..... 622.
645 ..... 623.
646 ..... 624.
647 ..... 625.
648 ..... 626.
649 ..... 627.
650 ..... 628.
651 ..... 629.
652 ..... 630.
653 ..... 631.
654 ..... 632.
655 ..... 633.
656 ..... 634.
657 ..... 635.
658 ..... 636.
659 ..... 637.
660 ..... 638.
661 ..... 639.
662 ..... 640.
663 ..... 641.
664 ..... 642.
665 ..... 643.
666 ..... 644.
667 ..... 645.
668 ..... 646.
669 ..... 647.
670 ..... 648.
671 ..... 649.
672 ..... 650.
673 ..... 651.
674 ..... 652.
675 ..... 653.
676 ..... 654.
677 ..... 655.
678 ..... 656.
679 ..... 657.
680 ..... 658.
681 ..... 659.
682 ..... 660.
683 ..... 661.
684 ..... 662.
685 ..... 663.
686 ..... 664.
687 ..... 665.
688 ..... 666.
689 ..... 667.
690 ..... 668.
691 ..... 669.
692 ..... 670.
693 ..... 671.
694 ..... 672.
695 ..... 673.
696 ..... 674.
697 ..... 675.
698 ..... 676.
699 ..... 677.
700 ..... 678.
701 ..... 679.
702 ..... 680.
703 ..... 681.
704 ..... 682.
705 ..... 683.
706 ..... 684.
707 ..... 685.
708 ..... 686.
709 ..... 687.
710 ..... 688.
711 ..... 689.
712 ..... 690.
713 ..... 691.
714 ..... 692.
715 ..... 693.
716 ..... 694.
717 ..... 695.
718 ..... 696.
719 ..... 697.
720 ..... 698.
721 ..... 699.
722 ..... 700.
723 ..... 701.
724 ..... 702.
725 ..... 703.
726 ..... 704.
727 ..... 705.
728 ..... 706.
729 ..... 707.
730 ..... 708.
731 ..... 709.
732 ..... 710.
733 ..... 711.
734 ..... 712.
735 ..... 713.
736 ..... 714.
737 ..... 715.
738 ..... 716.
739 ..... 717.
740 ..... 718.
741 ..... 719.
742 ..... 720.
743 ..... 721.
744 ..... 722.
745 ..... 723.
746 ..... 724.
747 ..... 725.
748 ..... 726.
749 ..... 727.
750 ..... 728.
751 ..... 729.
752 ..... 730.
753 ..... 731.
754 ..... 732.
755 ..... 733.
756 ..... 734.
757 ..... 735.
758 ..... 736.
759 ..... 737.
760 ..... 738.
761 ..... 739.
762 ..... 740.
763 ..... 741.
764 ..... 742.
765 ..... 743.
766 ..... 744.
767 ..... 745.
768 ..... 746.
769 ..... 747.
770 ..... 748.
771 ..... 749.
772 ..... 750.
773 ..... 751.
774 ..... 752.
775 ..... 753.
776 ..... 754.
777 ..... 755.
778 ..... 756.
779 ..... 757.
780 ..... 758.
781 ..... 759.
782 ..... 760.
783 ..... 761.
784 ..... 762.
785 ..... 763.
786 ..... 764.
787 ..... 765.
788 ..... 766.
789 ..... 767.
790 ..... 768.
791 ..... 769.
792 ..... 770.
793 ..... 771.
794 ..... 772.
795 ..... 773.
796 ..... 774.
797 ..... 775.
798 ..... 776.
799 ..... 777.
800 ..... 778.
801 ..... 779.
802 ..... 780.
803 ..... 781.
804 ..... 782.
805 ..... 783.
806 ..... 784.
807 ..... 785.
808 ..... 786.
809 ..... 787.
810 ..... 788.
811 ..... 789.
812 ..... 790.
813 ..... 791.
814 ..... 792.
815 ..... 793.
816 ..... 794.
817 ..... 795.
818 ..... 796.
819 ..... 797.
820 ..... 798.
821 ..... 799.
822 ..... 800.
823 ..... 801.
824 ..... 802.
825 ..... 803.
826 ..... 804.
827 ..... 805.
828 ..... 806.
829 ..... 807.
830 ..... 808.
831 ..... 809.
832 ..... 810.
833 ..... 811.
834 ..... 812.
835 ..... 813.
836 ..... 814.
837 ..... 815.
838 ..... 816.
839 ..... 817.
840 ..... 818.
841 ..... 819.
842 ..... 820.
843 ..... 821.
844 ..... 822.
845 ..... 823.
846 ..... 824.
847 ..... 825.
848 ..... 826.
849 ..... 827.
850 ..... 828.
851 ..... 829.
852 ..... 830.
853 ..... 831.
854 ..... 832.
855 ..... 833.
856 ..... 834.
857 ..... 835.
858 ..... 836.
859 ..... 837.
860 ..... 838.
861 ..... 839.
862 ..... 840.
863 ..... 841.
864 ..... 842.
865 ..... 843.
866 ..... 844.
867 ..... 845.
868 ..... 846.
869 ..... 847.
870 ..... 848.
871 ..... 849.
872 ..... 850.
873 ..... 851.
874 ..... 852.
875 ..... 853.
876 ..... 854.
877 ..... 855.
878 ..... 856.
879 ..... 857.
880 ..... 858.
881 ..... 859.
882 ..... 860.
883 ..... 861.
884 ..... 862.
885 ..... 863.
886 ..... 864.
887 ..... 865.
888 ..... 866.
889 ..... 867.
890 ..... 868.
891 ..... 869.
892 ..... 870.
893 ..... 871.
894 ..... 872.
895 ..... 873.
896 ..... 874.
897 ..... 875.
898 ..... 876.
899 ..... 877.
900 ..... 878.
901 ..... 879.
902 ..... 880.
903 ..... 881.
904 ..... 882.
905 ..... 883.
906 ..... 884.
907 ..... 885.
908 ..... 886.
909 ..... 887.
910 ..... 888.
911 ..... 889.
912 ..... 890.
913 ..... 891.
914 ..... 892.
915 ..... 893.
916 ..... 894.
917 ..... 895.
918 ..... 896.
919 ..... 897.
920 ..... 898.
921 ..... 899.
922 ..... 900.
923 ..... 901.
924 ..... 902.
925 ..... 903.
926 ..... 904.
927 ..... 905.
928 ..... 906.
929 ..... 907.
930 ..... 908.
931 ..... 909.
932 ..... 910.
933 ..... 911.
934 ..... 912.
935 ..... 913.
936 ..... 914.
937 ..... 915.
938 ..... 916.
939 ..... 917.
940 ..... 918.
941 ..... 919.
942 ..... 920.
943 ..... 921.
944 ..... 922.
945 ..... 923.
946 ..... 924.
947 ..... 925.
948 ..... 926.
949 ..... 927.
950 ..... 928.
951 ..... 929.
952 ..... 930.
953 ..... 931.
954 ..... 932.
955 ..... 933.
956 ..... 934.
957 ..... 935.
958 ..... 936.
959 ..... 937.
960 ..... 938.
961 ..... 939.
962 ..... 940.
963 ..... 941.
964 ..... 942.
965 ..... 943.
966 ..... 944.
967 ..... 945.
968 ..... 946.
969 ..... 947.
970 ..... 948.
971 ..... 949.
972 ..... 950.
973 ..... 951.
974 ..... 952.
975 ..... 953.
976 ..... 954.
977 ..... 955.
978 ..... 956.
979 ..... 957.
980 ..... 958.
981 ..... 959.
982 ..... 960.
983 ..... 961.
984 ..... 962.
985 ..... 963.
986 ..... 964.
987 ..... 965.
988 ..... 966.
989 ..... 967.
990 ..... 968.
991 ..... 969.
992 ..... 970.
993 ..... 971.
994 ..... 972.
995 ..... 973.
996 ..... 974.
997 ..... 975.
998 ..... 976.
999 ..... 977.
1000 ..... 978.
1001 ..... 979.
1002 ..... 980.
1003 ..... 981.
1004 ..... 982.
1005 ..... 983.
1006 ..... 984.
1007 ..... 985.
1008 ..... 986.
1009 ..... 987.
1010 ..... 988.
1011 ..... 989.
1012 ..... 990.
1013 ..... 991.
1014 ..... 992.
1015 ..... 993.
1016 ..... 994.
1017 ..... 995.
1018 ..... 996.
1019 ..... 997.
1020 ..... 998.
1021 ..... 999.
1022 ..... 1000.
1023 ..... 1001.
1024 ..... 1002.
1025 ..... 1003.
1026 ..... 1004.
1027 ..... 1005.
1028 ..... 1006.
1029 ..... 1007.
1030 ..... 1008.
1031 ..... 1009.
1032 ..... 1010.
1033 ..... 1011.
1034 ..... 1012.
1035 ..... 1013.
1036 ..... 1014.
1037 ..... 1015.
1038 ..... 1016.
1039 ..... 1017.
1040 ..... 1018.
1041 ..... 1019.
1042 ..... 1020.
1043 ..... 1021.
1044 ..... 1022.
1045 ..... 1023.
1046 ..... 1024.
1047 ..... 1025.
1048 ..... 1026.
1049 ..... 1027.
1050 ..... 1028.
1051 ..... 1029.
1052 ..... 1030.
1053 ..... 1031.
1054 ..... 1032.
1055 ..... 1033.
1056 ..... 1034.
1057 ..... 1035.
1058 ..... 1036.
1059 ..... 1037.
1060 ..... 1038.
1061 ..... 1039.
1062 ..... 1040.
1063 ..... 1041.
1064 ..... 1042.
1065 ..... 1043.
1066 ..... 1044.
1067 ..... 1045.
1068 ..... 1046.
1069 ..... 1047.
1070 ..... 1048.
1071 ..... 1049.
1072 ..... 1050.
1073 ..... 1051.
1074 ..... 1052.
1075 ..... 1053.
1076 ..... 1054.
1077 ..... 1055.
1078 ..... 1056.
1079 ..... 1057.
1080 ..... 1058.
1081 ..... 1059.
1082 ..... 1060.
1083 ..... 1061.
1084 ..... 1062.
1085 ..... 1063.
1086 ..... 1064.
1087 ..... 1065.
1088 ..... 1066.
1089 ..... 1067.
1090 ..... 1068.
1091 ..... 1069.
1092 ..... 1070.
1093 ..... 1071.
1094 ..... 1072.
1095 ..... 1073.
1096 ..... 1074.
1097 ..... 1075.
1098 ..... 1076.
1099 ..... 1077.
1100 ..... 1078.
1101 ..... 1079.
1102 ..... 1080.
1103 ..... 1081.
1104 ..... 1082.
1105 ..... 1083.
1106 ..... 1084.
1107 ..... 1085.
1108 ..... 1086.
1109 ..... 1087.
1110 ..... 1088.
1111 ..... 1089.
1112 ..... 1090.
1113 ..... 1091.
1114 ..... 1092.
1115 ..... 1093.
1116 ..... 1094.
1117 ..... 1095.
1118 ..... 1096.
1119 ..... 1097.
1120 ..... 1098.
1121 ..... 1099.
1122 ..... 1100.
1123 ..... 1101.
1124 ..... 1102.
1125 ..... 1103.
1126 ..... 1104.
1127 ..... 1105.
1128 ..... 1106.
1129 ..... 1107.
1130 ..... 1108.
1131 ..... 1109.
1132 ..... 1110.
1133 ..... 1111.
1134 ..... 1112.
1135 ..... 1113.
1136 ..... 1114.
1137 ..... 1115.
1138 ..... 1116.
1139 ..... 1117.
1140 ..... 1118.
1141 ..... 1119.
1142 ..... 1120.
1143 ..... 1121.
1144 ..... 1122.
1145 ..... 1123.
1146 ..... 1124.
1147 ..... 1125.
1148 ..... 1126.
1149 ..... 1127.
1150 ..... 1128.
1151 ..... 1129.
1152 ..... 1130.
1153 ..... 1131.
1154 ..... 1132.
1155 ..... 1133.
1156 ..... 1134.
1157 ..... 1135.
1158 ..... 1136.
1159 ..... 1137.
1160 ..... 1138.
1161 ..... 1139.
1162 ..... 1140.
1163 ..... 1141.
1164 ..... 1142.
1165 ..... 1143.
1166 ..... 1144.
1167 ..... 1145.
1168 ..... 1146.
1169 ..... 1147.
1170 ..... 1148.
1171 ..... 1149.
1172 ..... 1150.
1173 ..... 1151.
1174 ..... 1152.
1175 ..... 1153.
1176 ..... 1154.
1177 ..... 1155.
1178 ..... 1156.
1179 ..... 1157.
1180 ..... 1158.
1181 ..... 1159.
1182 ..... 1160.
1183 ..... 116
```

```

95 define TARGET_MEL_MIT 4 /* bit ratio of target relations */
96 define QUERY_ATT_MIT 4 /* bit ratio of queried attributes */
97 define TARGET_ATT_MIT 4 /* bit ratio of updated attributes */
98 define UPD_ATT_MIT 4
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188

```

```

189 define GBJ_SYN 4
190 define GBJ_SYN 4
191 define GBJ_SYN 4
192 define PP 7
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283

```



```

847 else { /* not a bus queue, pick the first one in queue */
848   RDBG (f) = QLINK (empl_k);
849   /* empl_k = index of entry dequeued */
850   /* DDD */
851   if (empl_tr > 0)
852     /* EEE */
853     /* EEE */
854     /* EEE */
855     /* EEE */
856     /* EEE */
857     /* EEE */
858     /* EEE */
859     /* EEE */
860     /* EEE */
861     /* EEE */
862     /* EEE */
863     /* EEE */
864     /* EEE */
865     /* EEE */
866     /* EEE */
867     /* EEE */
868     /* EEE */
869     /* EEE */
870     /* EEE */
871     /* EEE */
872     /* EEE */
873     /* EEE */
874     /* EEE */
875     /* EEE */
876     /* EEE */
877     /* EEE */
878     /* EEE */
879     /* EEE */
880     /* EEE */
881     /* EEE */
882     /* EEE */
883     /* EEE */
884     /* EEE */
885     /* EEE */
886     /* EEE */
887     /* EEE */
888     /* EEE */
889     /* EEE */
890     /* EEE */
891     /* EEE */
892     /* EEE */
893     /* EEE */
894     /* EEE */
895     /* EEE */
896     /* EEE */
897     /* EEE */
898     /* EEE */
899     /* EEE */
900     /* EEE */
901     /* EEE */
902     /* EEE */
903     /* EEE */
904     /* EEE */
905     /* EEE */
906     /* EEE */
907     /* EEE */
908     /* EEE */
909     /* EEE */
910     /* EEE */
911     /* EEE */
912     /* EEE */
913     /* EEE */
914     /* EEE */
915     /* EEE */
916     /* EEE */
917     /* EEE */
918     /* EEE */
919     /* EEE */
920     /* EEE */
921     /* EEE */
922     /* EEE */
923     /* EEE */
924     /* EEE */
925     /* EEE */
926     /* EEE */
927     /* EEE */
928     /* EEE */
929     /* EEE */
930     /* EEE */
931     /* EEE */
932     /* EEE */
933     /* EEE */
934     /* EEE */
935     /* EEE */
936     /* EEE */
937     /* EEE */
938     /* EEE */
939     /* EEE */
940     /* EEE */

```

```

941 RDBG (f) /* QLEN (f) * (clock - TABF (f)) /* queue state */
942 TABF (f) = clock;
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023

```

```

1024 bro_activate [tr] [site] [acc] = WAIT_BRO;
1025 for all pre
1026   if ( relation_count [pre] [site] == 1
1027     || relation_count [pre] [site] == 1
1028     || duration_count [tr] [site] [acc] == 1
1029     || duration_count [tr] [site] [acc] == 1
1030   )
1031     component_acc [tr] [site] [acc] = 1;
1032   component_acc [tr] [site] [acc] = 0;
1033   }
1034   }
1035   }
1036   }
1037   for_all_acc for_all_pre bro_component [tr] [acc] [pre] = NULL;
1038 }
1039
1040 /-----
1041 * Main simulation program : sim ()
1042 *
1043 *
1044 *
1045 *
1046 *
1047 *
1048 *
1049 *
1050 *
1051 *
1052 *
1053 *
1054 *
1055 *
1056 *
1057 *
1058 *
1059 *
1060 *
1061 *
1062 *
1063 *
1064 *
1065 *
1066 *
1067 *
1068 *
1069 *
1070 *
1071 *
1072 *
1073 *
1074 *
1075 *
1076 *
1077 *
1078 *
1079 *
1080 *
1081 *
1082 *
1083 *
1084 *
1085 *
1086 *
1087 *
1088 *
1089 *
1090 *
1091 *
1092 *
1093 *
1094 *
1095 *
1096 *
1097 *
1098 *
1099 *
1100 *
1101 *
1102 *
1103 *
1104 *
1105 *
1106 *
1107 *
1108 *
1109 *
1110 *
1111 *
1112 *
1113 *
1114 *
1115 *
1116 *
1117 *
1118 *
1119 *
1120 *
1121 *

```

```

1118 continue;
1119 }
1120 switch ( event ) {
1121 case ( NEW_TR ) :
1122   new_cram ();
1123   break;
1124 case ( NEED_BUS_TOKEN ) :
1125   if ( need_bus_token ( token ) == 0 ) return;
1126   if ( reserve ( bus_token [ site ], token, 0 ) == BUSY ) {
1127     result = 0;
1128     return;
1129   }
1130   break;
1131 case ( CONTEND_BUS ) :
1132   if ( reserve ( bus_token [ site ], token, 0 ) == BUSY ) {
1133     result = 0;
1134     return;
1135   }
1136   break;
1137 case ( RELEASE_BUS_TOKEN ) :
1138   release ( bus_token [ site ], token );
1139   break;
1140 case ( RELEASE_BRO ) :
1141   release ( bro_token [ site ], token );
1142   break;
1143 case ( RELEASE_DISK ) :
1144   release ( disk [ site ], token );
1145   break;
1146 case ( RELEASE_BUS ) :
1147   release ( bus_token [ site ], token );
1148   break;
1149 case ( NEED_DISK ) :
1150   if ( reserve ( disk [ site ], token, 0 ) == BUSY ) {
1151     result = 0;
1152     return;
1153   }
1154   break;
1155 case ( RELEASE_DISK ) :
1156   release ( disk [ site ], token );
1157   break;
1158 case ( RELEASE_DISK ) :
1159   release ( disk [ site ], token );
1160   break;
1161 case ( RELEASE_DISK ) :
1162   release ( disk [ site ], token );
1163   break;
1164 case ( RELEASE_DISK ) :
1165   release ( disk [ site ], token );
1166   break;
1167 case ( LP_FINISH ) :
1168   finish ( token );
1169   break;
1170 case ( GP_FINISH ) :
1171   finish ( token );
1172   break;
1173 case ( PR_FINISH ) :
1174   finish ( token );
1175   break;
1176 case ( WAKE_UP ) :
1177   wake_up ( token );
1178   break;
1179 case ( STOP ) :
1180   goto stop;
1181 }
1182 }
1183 }
1184 }
1185 }
1186 }
1187 }
1188 }
1189 }
1190 }
1191 }
1192 }
1193 }
1194 }
1195 }
1196 }
1197 }
1198 }
1199 }
1200 }
1201 }
1202 }
1203 }
1204 }
1205 }
1206 }
1207 }
1208 }
1209 }
1210 }
1211 }

```



```

1400 * count the cardinality. The different copies may reach a
1401 * different number of copies.
1402 * RI and PI Phase. That also make diff copies different.
1403 * Here we always use the first site that contains the
1404 * updated relation
1405
1406 { int bit_per_tuple, card, size, rsl;
1407
1408     size = 0;
1409     for_all_site [ tr ]
1410     for_all_site
1411     if ( rsl [ rsl ] == 1 ) break;
1412
1413     card = 1;
1414     for_all_site
1415     if ( updated_jact [ tr ] [ rsl ] == 1 )
1416         card = site_act.card [ tr ] [ site ] [ act ];
1417
1418     bit_per_tuple = 0;
1419     for_all_site
1420     if ( updated_jact [ tr ] [ rsl ] [ act ] == 1 )
1421         if ( bit_per_tuple * "down" >= bit_per_tuple )
1422             bit_per_tuple = "down" > bit_per_tuple;
1423     bit_per_tuple = relation_j_rsl [ rsl ];
1424
1425     packet = ( card * bit_per_tuple ) / DATA_SIZE;
1426     packet_max ( packet );
1427     schedule_B ( RELEASE_BUS, packet, "PACKET_TIME", token );
1428     break;
1429 }
1430
1431 /*----- after using bus
1432 -----*/
1433
1434 release_bus ();
1435
1436 release_bus ( token );
1437
1438 release_bus ( token );
1439
1440 register_int ( tr, rsl, site_act );
1441
1442 site = site_of [ token ];
1443
1444 tr = token;
1445
1446 switch ( site_act [ tr ] [ rsl ] ) {
1447
1448     /*----- LP phase
1449     -----*/
1450     case LP_WAIT_LOCK
1451         /* The following 3 state is for trans shoes
1452         * lock had been withdrwn
1453         * due to a PG conflict
1454         */
1455         case LP_WAIT_LOCK
1456             for_all_site [ blocked_site [ tr ] ]
1457             break;
1458         case LP_WAIT_LOCK
1459             for_all_site [ blocked_site [ tr ] ]
1460             break;
1461         case LP_WAIT_LOCK
1462             for_all_site [ blocked_site [ tr ] ]
1463             break;
1464         case LP_WAIT_LOCK
1465             for_all_site [ blocked_site [ tr ] ]
1466             break;
1467         case LP_WAIT_LOCK
1468             for_all_site [ blocked_site [ tr ] ]
1469             break;
1470         case LP_WAIT_LOCK
1471             for_all_site [ blocked_site [ tr ] ]
1472             break;
1473         case LP_WAIT_LOCK
1474             for_all_site [ blocked_site [ tr ] ]
1475             break;
1476         case LP_WAIT_LOCK
1477             for_all_site [ blocked_site [ tr ] ]
1478             break;
1479         case LP_WAIT_LOCK
1480             for_all_site [ blocked_site [ tr ] ]
1481             break;
1482         case LP_WAIT_LOCK
1483             for_all_site [ blocked_site [ tr ] ]
1484             break;
1485         case LP_WAIT_LOCK
1486             for_all_site [ blocked_site [ tr ] ]
1487             break;
1488         case LP_WAIT_LOCK
1489             for_all_site [ blocked_site [ tr ] ]
1490             break;
1491         case LP_WAIT_LOCK
1492             for_all_site [ blocked_site [ tr ] ]
1493             break;

```

```

1494 child_gen ( tr );
1495         break;
1496     }
1497     return;
1498 }
1499
1500 goto_site ( LP_INIT );
1501 num_in_ip [ tr ] ++;
1502 child_gen ( tr );
1503 pr ( "Home" lock );
1504
1505 case LP_INIT /* entry point of a new child */
1506     if ( need_lock_in_ip ( token ) == 0 ) {
1507         schedule_C ( NEED_DISK, 1000 * INSTRUCTION_TIME, token );
1508         return;
1509     }
1510     goto_site ( LP_WAIT_LOCK );
1511     wake_up ( token );
1512
1513 /*----- the only one child site home site is empty
1514 -----*/
1515 case LP_BJ
1516     /* lock ( tr, site );
1517     goto_site ( RI_START );
1518     last_pt = 0;
1519     goto_orphan, cvtln;
1520     break;
1521
1522 case LP_BYM_EST_RND
1523     pr ( ESTIMATE_SLOWEST_SITE Broadcast HW );
1524     cc_lock ( tr, site );
1525     goto_site ( RI_START );
1526     last_pt = 0;
1527     goto_orphan, cvtln;
1528     break;
1529
1530 /* need check if all have finished or not */
1531 if ( num_in_ip [ tr ] == 0 ) return;
1532 num_in_ip [ tr ] = "LP_WAIT_LOCK";
1533
1534 if ( num_in_ip [ tr ] == 0 ) return;
1535 num_of_broad [ tr ] = 1;
1536 goto_site ( LP_ALL_FINISH );
1537 blocked_site [ tr ] = 0;
1538
1539 case LP_ALL_FINISH
1540     pr ( LP_ALL_FINISH );
1541     schedule_C ( NEED_BUS, TOKEN, DELTA_TIME, token );
1542     break;
1543
1544 case LP_WAIT_LOCK
1545     pr ( Alternate ESTIMATE SLOWEST SITE Broadcast HW );
1546     cc_lock ( tr, site );
1547     if ( r_unlock_time == ANYTIME ) {
1548         cc_unlock ( tr, site );
1549         resume_by ( tr, site );
1550         break;
1551     }
1552     case LP_WAIT_LOCK
1553         /* start GSI phase */
1554         cc_lock ( tr, site );
1555         if ( r_unlock_time == ANYTIME ) {
1556             cc_unlock ( tr, site );
1557             resume_by ( tr, site );
1558             break;
1559         }
1560         goto_site ( GSI );
1561         goto_site ( LP_WAIT_LOCK );
1562         I = pick_min_site ( tr, site, LP_WAIT_LOCK );
1563         wake_up ( token_of [ tr ] [ I ] );
1564         break;
1565
1566 case GSI_INIT
1567     /* ONE GSI PHASE A MIN SITE IS IDENTIFIED */
1568     /*

```



```

1776         || is_non_update_site ( i )
1777     )
1778     terminate_site ( token_of [ tr ] [ i ] , tr , i ) ;
1779     else go_sleep ( tr , i , PP_WAIT_UPDATE ) ;
1780 }
1781 }
1782 if ( last_tr == 0 ) {
1783     goto_phase ( PP ) ;
1784     wake_up ( token_of [ tr ] [ home_of [ tr ] ] ) ;
1785 }
1786 /*
1787  * Wakeup home site for PP
1788 */
1789 else {
1790     for_all_related_join_state ( RT_WAIT_START )
1791     wake_up ( token_of [ tr ] [ i ] ) ;
1792 }
1793 }
1794 }
1795 }
1796 }
1797 }
1798 }
1799 }
1800 }
1801 }
1802 }
1803 }
1804 }
1805 }
1806 }
1807 }
1808 }
1809 }
1810 }
1811 }
1812 }
1813 }
1814 }
1815 }
1816 }
1817 }
1818 }
1819 }
1820 }
1821 }
1822 }
1823 }
1824 }
1825 }
1826 }
1827 }
1828 }
1829 }
1830 }
1831 }
1832 }
1833 }
1834 }
1835 }
1836 }
1837 }
1838 }
1839 }
1840 }
1841 }
1842 }
1843 }
1844 }
1845 }
1846 }
1847 }
1848 }
1849 }
1850 }
1851 }
1852 }
1853 }
1854 }
1855 }
1856 }
1857 }
1858 }
1859 }
1860 }
1861 }
1862 }
1863 }
1864 }
1865 }
1866 }
1867 }
1868 }
1869 }
1870 }
1871 }
1872 }
1873 }
1874 }
1875 }
1876 }
1877 }
1878 }
1879 }
1880 }
1881 }
1882 }
1883 }
1884 }
1885 }
1886 }
1887 }
1888 }
1889 }
1890 }
1891 }
1892 }
1893 }
1894 }
1895 }
1896 }
1897 }
1898 }
1899 }
1900 }
1901 }
1902 }
1903 }
1904 }
1905 }
1906 }
1907 }
1908 }
1909 }
1910 }
1911 }
1912 }
1913 }
1914 }
1915 }
1916 }
1917 }
1918 }
1919 }
1920 }
1921 }
1922 }
1923 }
1924 }
1925 }
1926 }
1927 }
1928 }
1929 }
1930 }
1931 }
1932 }
1933 }
1934 }
1935 }
1936 }
1937 }
1938 }
1939 }
1940 }
1941 }
1942 }
1943 }
1944 }
1945 }
1946 }
1947 }
1948 }
1949 }
1950 }
1951 }
1952 }
1953 }
1954 }
1955 }
1956 }
1957 }
1958 }
1959 }
1960 }
1961 }
1962 }
1963 }

```

```

1870 }
1871 }
1872 }
1873 }
1874 }
1875 }
1876 }
1877 }
1878 }
1879 }
1880 }
1881 }
1882 }
1883 }
1884 }
1885 }
1886 }
1887 }
1888 }
1889 }
1890 }
1891 }
1892 }
1893 }
1894 }
1895 }
1896 }
1897 }
1898 }
1899 }
1900 }
1901 }
1902 }
1903 }
1904 }
1905 }
1906 }
1907 }
1908 }
1909 }
1910 }
1911 }
1912 }
1913 }
1914 }
1915 }
1916 }
1917 }
1918 }
1919 }
1920 }
1921 }
1922 }
1923 }
1924 }
1925 }
1926 }
1927 }
1928 }
1929 }
1930 }
1931 }
1932 }
1933 }
1934 }
1935 }
1936 }
1937 }
1938 }
1939 }
1940 }
1941 }
1942 }
1943 }
1944 }
1945 }
1946 }
1947 }
1948 }
1949 }
1950 }
1951 }
1952 }
1953 }
1954 }
1955 }
1956 }
1957 }
1958 }
1959 }
1960 }
1961 }
1962 }
1963 }

```



```

2329         if ( !is_read_only ) {
2330             terminate_site ( token, tr, site );
2331             terminate_tr ( tr );
2332             return;
2333         }
2334     }
2335     case PP_PRE_LOCK : /* entry point of waken process */
2336         goto_state ( PP_PRE_LOCK );
2337     }
2338     case PP_LOCK : /* entry point of waken process */
2339         switch ( get_local_lock ( token, tr, site ) ) {
2340             case REJECT :
2341                 go_sleep ( tr, site, PP_WAIT_LOCK );
2342                 prs ( BLOCKED ( tr ) );
2343                 blocked_site ( tr );
2344                 break;
2345             case GOT_LOCK :
2346                 if ( class_of [ tr ] != REMOTE ) {
2347                     int block_someone = 0;
2348                     block_someone = 0;
2349                     repeat ( TR )
2350                         if ( pg [ tr ] [ i ] { tr } == i )
2351                             || ( get_relative [ i ] { tr } == -1 )
2352                                 ) {
2353                         block_someone = 1;
2354                         break;
2355                     }
2356                 }
2357                 if ( block_someone == 1 ) /* block someone */
2358                     goto broad;
2359                 else {
2360                     /* Did you block any body? lock up lock
2361                     * NO user and read blocking message
2362                     * cpu time is ignored
2363                     * notice no disk read is needed
2364                     * blocked_lock can protect the read set to be
2365                     * allowed
2366                     * the following check is a redundancy checking
2367                     * since pg someone should hold lock
2368                     * the same function code
2369                     * the code is questionable No test yet
2370                     */
2371                     if ( cc_lock_get, but ( tr, site ) == FAIL ) {
2372                         double P;
2373                         P = 20 * INSTRUCTION_TIME * cc_time;
2374                         P = 20 * ( REVERSE_LSI * token );
2375                         return;
2376                     }
2377                     cc_lock ( tr, site );
2378                     cc_unlock ( tr, site );
2379                     resume_by ( tr, site );
2380                     terminate_site ( token, tr, site );
2381                     break;
2382                 }
2383                 /* Local query */
2384                 schedule ( NEED_BUS_TOKEN, DELTA_TIME, token );
2385                 break;
2386             case HAVE :
2387                 case STATE :
2388                     terminate_tr ( tr );
2389                     break;
2390             }
2391         }
2392     }
2393     }
2394     }
2395     }
2396     }
2397     }
2398     }
2399     }
2400     }
2401     }
2402     }
2403     }
2404     }
2405     }
2406     }
2407     }
2408     }
2409     }
2410     }
2411     }
2412     }
2413     }
2414     }
2415     }
2416     }
2417     }
2418     }
2419     }
2420     }
2421     }
2422     }

```

```

2423     }
2424     }
2425     }
2426     }
2427     }
2428     }
2429     }
2430     }
2431     }
2432     }
2433     }
2434     }
2435     }
2436     }
2437     }
2438     }
2439     }
2440     }
2441     }
2442     }
2443     }
2444     }
2445     }
2446     }
2447     }
2448     }
2449     }
2450     }
2451     }
2452     }
2453     }
2454     }
2455     }
2456     }
2457     }
2458     }
2459     }
2460     }
2461     }
2462     }
2463     }
2464     }
2465     }
2466     }
2467     }
2468     }
2469     }
2470     }
2471     }
2472     }
2473     }
2474     }
2475     }
2476     }
2477     }
2478     }
2479     }
2480     }
2481     }
2482     }
2483     }
2484     }
2485     }
2486     }
2487     }
2488     }
2489     }
2490     }
2491     }
2492     }
2493     }
2494     }
2495     }
2496     }
2497     }
2498     }
2499     }
2500     }
2501     }
2502     }
2503     }
2504     }
2505     }
2506     }
2507     }
2508     }
2509     }
2510     }
2511     }
2512     }
2513     }
2514     }
2515     }
2516     }

```

```

2617 * otherwise stay at RT_WAIT_START state
2618 */
2619
2620 if ( site == home_of [ tr ] ) {
2621   go_sleep ( tr, site, PP_WAIT_START );
2622   /* pseudo broadcast: the relations */
2623   for_all_rel [ rel ] in site_rel_queries {
2624     if ( ! rel [ tr ] [ rel ] == HAD_BRO )
2625       return;
2626   }
2627   if ( first_rel_site [ tr ] == 0 )
2628     first_rel [ tr ];
2629   if ( ! site == first_rel_site [ tr ] ) {
2630     goto_sleep ( tr, site, RT_WAIT_START );
2631     return;
2632   }
2633   goto_state ( RT );
2634   schedule_C ( NEED_DISK, DELTA_TIME, token );
2635 }
2636
2637 case RT_WAIT_START /* non home site only */
2638 :
2639   goto_sleep ( RT );
2640   schedule_C ( NEED_DISK, DELTA_TIME, token );
2641 }
2642
2643 /*-----*/
2644 /*-----*/
2645 /*-----*/
2646 /*-----*/
2647 /*-----*/
2648 /*-----*/
2649 /*-----*/
2650 /*-----*/
2651 /*-----*/
2652 /*-----*/
2653 /*-----*/
2654 /*-----*/
2655 /*-----*/
2656 /*-----*/
2657 /*-----*/
2658 /*-----*/
2659 /*-----*/
2660 /*-----*/
2661 /*-----*/
2662 /*-----*/
2663 /*-----*/
2664 /*-----*/
2665 /*-----*/
2666 /*-----*/
2667 /*-----*/
2668 /*-----*/
2669 /*-----*/
2670 /*-----*/
2671 /*-----*/
2672 /*-----*/
2673 /*-----*/
2674 /*-----*/
2675 /*-----*/
2676 /*-----*/
2677 /*-----*/
2678 /*-----*/
2679 /*-----*/
2680 /*-----*/
2681 /*-----*/
2682 /*-----*/
2683 /*-----*/
2684 /*-----*/
2685 /*-----*/
2686 /*-----*/
2687 /*-----*/
2688 /*-----*/
2689 /*-----*/
2690 /*-----*/
2691 /*-----*/
2692 /*-----*/
2693 /*-----*/
2694 /*-----*/
2695 /*-----*/
2696 /*-----*/
2697 /*-----*/
2698 /*-----*/
2699 /*-----*/
2700 /*-----*/
2701 /*-----*/
2702 /*-----*/
2703 /*-----*/
2704 /*-----*/

```

```

2611 repeat ( TR )
2612   old_pg [ 1 ] = pg [ tr ] [ 1 ];
2613   loop ( 1, 1, TR )
2614     loop ( 1, 1, TR )
2615       pg [ 1 ] [ 1 ] = new_pg [ 1 ] [ 1 ];
2616
2617 /*
2618  * Invalidated those temporary r-locks which is conflicting
2619  * with locks used by this trans and is hold by a trans
2620  * which had no precedence relationship with current trans before
2621  * but is preceded by it now
2622  * Don't care about directive precedence relationship
2623  * since no conflict may arise between them
2624  * This should be moved to the time after the bus is released
2625  */
2626 loop ( 1, 1, SITE ) /* check site by site to save time */
2627 :
2628   if ( 1 == site )
2629     pguser [ 1 ]
2630     pg [ tr ] [ pguser [ 1 ] ] == theformer
2631     old_pg [ pguser [ 1 ] ] == 0
2632     lock_validity [ pguser [ 1 ] ] [ 1 ] == VALID
2633   ) {
2634     /* found a trans in site 1 that is easily preceded by tr */
2635     trj = pguser [ 1 ];
2636     /* to see if the locks obtained in this site by trj is invalid */
2637     for_all_rel if ( queried_rel [ tr ] [ rel ] == 1 )
2638       if ( ( ( lock_state [ tr ] [ rel ] == READING
2639             || lock_state [ tr ] [ rel ] == READING_WAIT )
2640             || local_lock [ 1 ] [ rel ] == READING_WAIT )
2641             || local_lock [ 1 ] [ rel ] == READING_WAIT )
2642           ) {
2643             lock_state [ tr ] [ rel ] == WRITING
2644             || lock_state [ tr ] [ rel ] == READING_WAIT
2645             || local_lock [ 1 ] [ rel ] == READING_WAIT
2646             || local_lock [ 1 ] [ rel ] == READING_WAIT )
2647           {
2648             lock_validity [ trj ] [ 1 ] = INVALID;
2649             break;
2650           }
2651         }
2652       }
2653     }
2654     }
2655     }
2656     }
2657     }
2658     }
2659     }
2660     }
2661     }
2662     }
2663     }
2664     }
2665     }
2666     }
2667     }
2668     }
2669     }
2670     }
2671     }
2672     }
2673     }
2674     }
2675     }
2676     }
2677     }
2678     }
2679     }
2680     }
2681     }
2682     }
2683     }
2684     }
2685     }
2686     }
2687     }
2688     }
2689     }
2690     }
2691     }
2692     }
2693     }
2694     }
2695     }
2696     }
2697     }
2698     }
2699     }
2700     }
2701     }
2702     }
2703     }
2704     }

```

```

2706 * Save PG info in order to be restored at the end
2707 * If PG is not in the cache, this moment
2708 * It can only be stored into pg after release bus
2709 * This overhead is due to that cc_lock_get_bus uses pg not new_pg
2710
2711 loop ( 1, 1, TR )
2712   new_pg [ 1 ] [ 1 ] = PG [ 1 ] [ 1 ]
2713   repeat ( TR )
2714     old_pg [ 1 ] = PG [ tr ] [ 1 ]
2715     loop ( j, 1, TR ) {
2716       if ( new_relative [ tr ] [ j ] == 1 ) {
2717         if ( PG [ tr ] [ j ] == -1 ) {
2718           new_relative [ tr ] [ j ]
2719         }
2720         if ( not_lock_in_lp [ token_of [ tr ] [ site ] ] == 1 )
2721           ch_state [ tr ] [ site ] = LP_WAITEND
2722         solve_pg_conflict ( tr, site ),
2723         cc_time = timestamp () - enter_time,
2724         return ( FAIL )
2725       }
2726       if ( new_relative [ tr ] [ j ] == -1 ) {
2727         if ( PG [ tr ] [ j ] == -1 ) {
2728           new_relative [ tr ] [ j ]
2729         }
2730         if ( not_lock_in_lp [ token_of [ tr ] [ site ] ] == 1 )
2731           ch_state [ tr ] [ site ] = LP_WAITEND
2732         solve_pg_conflict ( tr, site ),
2733         cc_time = timestamp () - enter_time,
2734         return ( FAIL )
2735       }
2736       if ( new_relative [ tr ] [ j ] == 1 ) {
2737         if ( PG [ tr ] [ j ] == -1 ) {
2738           new_relative [ tr ] [ j ]
2739         }
2740         if ( not_lock_in_lp [ token_of [ tr ] [ site ] ] == 1 )
2741           ch_state [ tr ] [ site ] = LP_WAITEND
2742         solve_pg_conflict ( tr, site ),
2743         cc_time = timestamp () - enter_time,
2744         return ( FAIL )
2745       }
2746     }
2747   }
2748   new_relative [ tr ] [ 1 ] = -UNED,
2749   }
2750 }
2751
2752 /*
2753 * update lock state, powner, be careful, may be called more than once
2754 * in every phase including lp phase
2755 */
2756 struct ( tr, phase [ tr ] ) {
2757   case lp /* Confirm locks for those sites had obtained locks */
2758     for all related i {
2759       if ( powner [ i ] [ 1 ] == tr
2760         && lock_validity [ tr ] [ i ] == VALID )
2761         pre ( READING )
2762         lock_validity [ tr ] [ i ] = FIRMED
2763         powner [ i ] = NULL
2764         for all rel {
2765           lock_state [ rel ] [ i ] = local_lock [ i ] [ rel ]
2766           lock_lock [ i ] [ rel ]
2767           local_lock [ i ] [ rel ] = NULL
2768         }
2769       }
2770     }
2771     break
2772   }
2773   case OSJ
2774     pre ( READING )
2775     powner [ site ] = NULL
2776     for all rel local_lock [ site ] [ rel ] = NULL
2777     local_resume ( site )
2778     break
2779   }
2780 }
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3830
3831
3832
3833
3834
3835
3836
3837
3838
3839
3840
3841
3842
3843
3844
3845
3846
3847
3848
3849
3850
3851
3852
3853
3854
3855
3856
3857
3858
3859
3860
3861
3862
3863
3864
3865
3866
3867
3868
3869
3870
3871
3872
3873
3874
3875
3876
3877
3878
3879
3880
3881
3882
3883
3884
3885
3886
3887
3888
3889
3890
3891
3892
3893
3894
3895
3896
3897
3898
3899
3900
3901
3902
3903
3904
3905
3906
3907
3908
3909
3910
3911
3912
3913
3914
3915
3916
3917
3918
3919
3920
3921
3922
3923
3924
3925
3926
3927
3928
3929
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949
3950
3951
3952
3953
3954
3955
3956
3957
3958
3959
3960
3961
3962
3963
3964
3965
3966
3967
3968
3969
3970
3971
3972
3973
3974
3975
3976
3977
3978
3979
3980
3981
3982
3983
3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3999
4000

```



```

2893 * Pp phase: unlock all r_locks and w_locks
2894 * In each phase, the set of locked relations are identified first
2895 * since a copy of the lock is copied into a second
2896 * of semi_join broadcast
2897 */
2898
2899 /*
2900 * r_lock: int tr, site
2901 */
2902 register int r_lock;
2903 register int w_lock;
2904 int to_be_unlock [ R + 1 ] ;
2905
2906 switch ( tr_phase [ tr ] ) {
2907   case LP
2908     {
2909       if ( r_lock_time == ANYTIME ) {
2910         /*
2911          * be careful, there may be some sites haven't got lock yet
2912          * there may be some sites haven't finish yet if called late
2913          * time when called later, careful too.
2914          * save all the lock status into the vector to be unlock
2915          * since those locks still used by some sites
2916          * release those locks no body use it any more
2917          */
2918         for_all_rel
2919           to_be_unlock [ rel ] = lock_state [ tr ] [ rel ] ;
2920
2921         for_all_related_i
2922           if ( the_state [ tr ] [ i ] == LP_PRE_LOCK
2923             || the_state [ tr ] [ i ] == LP_DISK_LOCK
2924             || the_state [ tr ] [ i ] == LP_DISK
2925           )
2926             for_all_rel
2927               if ( queried_rel [ tr ] [ rel ] == 1
2928                 && r_lock [ i ] == 1
2929                 && to_be_unlock [ rel ] == 0
2930               )
2931                 break;
2932             }
2933         }
2934
2935         case GR1
2936           if ( r_lock_time == ANYTIME ) {
2937             /*
2938              * Check if the relations locked by this site
2939              * are not in the same phase
2940              * if not, unlock them
2941              */
2942             for_all_rel
2943               to_be_unlock [ rel ] = lock_state [ tr ] [ rel ] ;
2944
2945             /*
2946              * The following piece of codes is redundant since once
2947              * not lock in the same phase, the relation will not set here
2948              * for the current implementation
2949              * (only one copy in the system)
2950              */
2951             for_all_related_i
2952               if ( the_state [ tr ] [ i ] == 1
2953                 && not lock_in_the_token_of [ tr ] [ i ] == 1
2954               )
2955                 for_all_rel
2956                   if ( queried_rel [ tr ] [ rel ] == 1
2957                     && r_lock [ i ] == 1
2958                     && to_be_unlock [ rel ] == 0
2959                   )
2960                     for_all_queried_rel_in_site
2961                       if ( relation_conf [ rel ] [ broadcast_all [ tr ] ] == 1
2962                         && to_be_unlock [ rel ] == READING
2963                       )
2964

```

```

2965
2966 /*
2967 * The r-locked rel contains the bro att, could be
2968 * a candidate for unlocking
2969 */
2970
2971 K = 0; /* check to see if all att had been bro */
2972
2973 for_all_att
2974   if ( relation_conf [ rel ] [ att ] == 1
2975     && queried_rel [ tr ] [ rel ] == 1
2976   )
2977     /* check if the att had been broadcast */
2978     if ( bro_component [ tr ] [ att ] [ rel ] == WAIT_BRO )
2979       K = 1;
2980   break;
2981 }
2982
2983 if ( K == 0 ) lock_state [ tr ] [ rel ] = 1;
2984
2985 }
2986 /* had r-locked */
2987 /* ANYTIME */
2988 if ( r_lock_time == OBJ_END ) {
2989   /*
2990    * consider those non-target relations only
2991    * check if all the att had been broadcast
2992    * if yes, unlock them
2993    */
2994   for_all_queried_rel_in_site
2995     if ( target_rel [ tr ] [ rel ] == 0 ) {
2996       if ( relation_conf [ rel ] [ broadcast_all [ tr ] ] == 1
2997         && lock_state [ tr ] [ rel ] == READING
2998       )
2999         /*
3000          * The r-locked rel contains the bro att, could be a
3001          * candidate for unlocking
3002          */
3003         K = 0; /* check to see if all att had been bro */
3004         for_all_att
3005           if ( relation_conf [ rel ] [ att ] == 1
3006             && queried_rel [ tr ] [ rel ] == 1
3007             && check_if_the_att_had_been_bro /*
3008               if ( bro_component [ tr ] [ att ] [ rel ] == WAIT_BRO )
3009                 K = 1;
3010             )
3011               break;
3012         }
3013         /*
3014          * The r-locked rel contains the bro att, could be a
3015          * candidate for unlocking
3016          */
3017         K = 0; /* check to see if all att had been bro */
3018         for_all_att
3019           if ( relation_conf [ rel ] [ att ] == 1
3020             && queried_rel [ tr ] [ rel ] == 1
3021             && check_if_the_att_had_been_bro /*
3022               if ( bro_component [ tr ] [ att ] [ rel ] == WAIT_BRO )
3023                 K = 1;
3024             )
3025               break;
3026         }
3027         if ( K == 0 ) lock_state [ tr ] [ rel ] = 1;
3028       }
3029     }
3030
3031     case RT
3032       for_all_rel
3033         if ( lock_state [ tr ] [ rel ] == READING
3034           && bro_state [ tr ] [ rel ] == BRO_BRO
3035         )
3036           break;
3037         lock_state [ tr ] [ rel ] = 1;
3038       }
3039     }
3040
3041     case PP
3042       /*
3043        * From now on, this transaction has nothing to do with
3044        * CC any more. It can be created as having been terminated
3045        */
3046

```

```

3070 */
3071 if (class_of [ cr ] == LOCAL_QUERY) {
3072     int l_lock;
3073     l_lock = lock_state [ cr ] [ writing_rel [ cr ] ];
3074     for_all_rel { tr } { rel } = NULL;
3075     lock_state [ cr ] [ writing_rel [ cr ] ] = l_lock;
3076 }
3077
3078 /*
3079  * child_gen () --- generate child processes
3080  * set up necessary information
3081  * CALLED BY: release_bur PRE
3082  */
3083
3084 pre ( cc_unlock );
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3830
3831
3832
3833
3834
3835
3836
3837
3838
3839
3840
3841
3842
3843
3844
3845
3846
3847
3848
3849
3850
3851
3852
3853
3854
3855
3856
3857
3858
3859
3860
3861
3862
3863
3864
3865
3866
3867
3868
3869
3870
3871
3872
3873
3874
3875
3876
3877
3878
3879
3880
3881
3882
3883
3884
3885
3886
3887
3888
3889
3890
3891
3892
3893
3894
3895
3896
3897
3898
3899
3900
3901
3902
3903
3904
3905
3906
3907
3908
3909
3910
3911
3912
3913
3914
3915
3916
3917
3918
3919
3920
3921
3922
3923
3924
3925
3926
3927
3928
3929
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949
3950
3951
3952
3953
3954
3955
3956
3957
3958
3959
3960
3961
3962
3963
3964
3965
3966
3967
3968
3969
3970
3971
3972
3973
3974
3975
3976
3977
3978
3979
3980
3981
3982
3983
3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3999
4000

```



```

3634         if ( (
3635             { db_state [tr] [site] != home_of [tr]
3636             || flag == 0
3637             || site == slowest_site
3638             } ) ) {
3639             J = 0;
3640             for_all_queried_rel_in_site
3641                 * J = 0;
3642             }
3643             if ( J > max_time ) {
3644                 * J = 0;
3645                 slowest_site = site;
3646             }
3647             if ( db_state [tr] [slowest_site] == LP_WAIT_END ) {
3648                 flag = 1;
3649                 /* avoid select a site already in LP_WAIT_END */
3650                 goto avoid;
3651             }
3652             get_slowest_site_of [tr] = slowest_site;
3653             }
3654             }
3655             }
3656             }
3657             }
3658             }
3659             }
3660             }
3661             }
3662             }
3663             }
3664             }
3665             }
3666             }
3667             }
3668             }
3669             }
3670             }
3671             }
3672             }
3673             }
3674             }
3675             }
3676             }
3677             }
3678             }
3679             }
3680             }
3681             }
3682             }
3683             }
3684             }
3685             }
3686             }
3687             }
3688             }
3689             }
3690             }
3691             }
3692             }
3693             }
3694             }
3695             }
3696             }
3697             }
3698             }
3699             }
3700             }
3701             }
3702             }
3703             }
3704             }
3705             }
3706             }
3707             }
3708             }
3709             }
3710             }
3711             }
3712             }
3713             }
3714             }
3715             }
3716             }
3717             }
3718             }
3719             }
3720             }
3721             }
3722             }
3723             }
3724             }
3725             }
3726             }
3727             }
3728             }
3729             }
3730             }
3731             }
3732             }
3733             }
3734             }
3735             }
3736             }
3737             }
3738             }
3739             }
3740             }
3741             }
3742             }
3743             }
3744             }
3745             }
3746             }
3747             }
3748             }
3749             }
3750             }
3751             }
3752             }
3753             }
3754             }
3755             }
3756             }
3757             }
3758             }
3759             }
3760             }
3761             }
3762             }
3763             }
3764             }
3765             }
3766             }
3767             }
3768             }
3769             }
3770             }
3771             }
3772             }
3773             }
3774             }
3775             }
3776             }
3777             }
3778             }
3779             }
3780             }
3781             }
3782             }
3783             }
3784             }
3785             }
3786             }
3787             }
3788             }
3789             }
3790             }
3791             }
3792             }
3793             }
3794             }
3795             }
3796             }
3797             }
3798             }
3799             }
3800             }
3801             }
3802             }
3803             }
3804             }
3805             }
3806             }
3807             }
3808             }
3809             }
3810             }
3811             }
3812             }
3813             }
3814             }
3815             }
3816             }
3817             }
3818             }
3819             }
3820             }
3821             }

```

```

3822 * get_token_id() --- Get a token from the token pool
3823 *
3824 *
3825 *
3826 *
3827 *
3828 *
3829 *
3830 *
3831 *
3832 *
3833 *
3834 *
3835 *
3836 *
3837 *
3838 *
3839 *
3840 *
3841 *
3842 *
3843 *
3844 *
3845 *
3846 *
3847 *
3848 *
3849 *
3850 *
3851 *
3852 *
3853 *
3854 *
3855 *
3856 *
3857 *
3858 *
3859 *
3860 *
3861 *
3862 *
3863 *
3864 *
3865 *
3866 *
3867 *
3868 *
3869 *
3870 *
3871 *
3872 *
3873 *
3874 *
3875 *
3876 *
3877 *
3878 *
3879 *
3880 *
3881 *
3882 *
3883 *
3884 *
3885 *
3886 *
3887 *
3888 *
3889 *
3890 *
3891 *
3892 *
3893 *
3894 *
3895 *
3896 *
3897 *
3898 *
3899 *
3900 *
3901 *
3902 *
3903 *
3904 *
3905 *
3906 *
3907 *
3908 *
3909 *
3910 *
3911 *
3912 *
3913 *
3914 *
3915 *

```

```

3916 *
3917 *
3918 *
3919 *
3920 *
3921 *
3922 *
3923 *
3924 *
3925 *
3926 *
3927 *
3928 *
3929 *
3930 *
3931 *
3932 *
3933 *
3934 *
3935 *
3936 *
3937 *
3938 *
3939 *
3940 *
3941 *
3942 *
3943 *
3944 *
3945 *
3946 *
3947 *
3948 *
3949 *
3950 *
3951 *
3952 *
3953 *
3954 *
3955 *
3956 *
3957 *
3958 *
3959 *
3960 *
3961 *
3962 *
3963 *
3964 *
3965 *
3966 *
3967 *
3968 *
3969 *
3970 *
3971 *
3972 *
3973 *
3974 *
3975 *
3976 *
3977 *
3978 *
3979 *
3980 *
3981 *
3982 *
3983 *
3984 *
3985 *
3986 *
3987 *
3988 *
3989 *
3990 *
3991 *
3992 *
3993 *
3994 *
3995 *
3996 *
3997 *
3998 *
3999 *
4000 *
4001 *
4002 *
4003 *
4004 *
4005 *
4006 *
4007 *
4008 *
4009 *

```

```

4010 for_all_att {
4011   if (relation_conf [rel] [att] == 0) {
4012     int rel_card [rel] = (domain_range [att]/K);
4013     }
4014   }
4015   }
4016   }
4017   }
4018   }
4019   }
4020   }
4021   }
4022   }
4023   }
4024   }
4025   }
4026   }
4027   }
4028   }
4029   }
4030   }
4031   }
4032   }
4033   }
4034   }
4035   }
4036   }
4037   }
4038   }
4039   }
4040   }
4041   }
4042   }
4043   }
4044   }
4045   }
4046   }
4047   }
4048   }
4049   }
4050   }
4051   }
4052   }
4053   }
4054   }
4055   }
4056   }
4057   }
4058   }
4059   }
4060   }
4061   }
4062   }
4063   }
4064   }
4065   }
4066   }
4067   }
4068   }
4069   }
4070   }
4071   }
4072   }
4073   }
4074   }
4075   }
4076   }
4077   }
4078   }
4079   }
4080   }
4081   }
4082   }
4083   }
4084   }
4085   }
4086   }
4087   }
4088   }
4089   }
4090   }
4091   }
4092   }

```

```

4093   }
4094   }
4095   }
4096   }
4097   }
4098   }
4099   }
4100   }
4101   }
4102   }
4103   }
4104   }
4105   }
4106   }
4107   }
4108   }
4109   }
4110   }
4111   }
4112   }
4113   }
4114   }
4115   }
4116   }
4117   }
4118   }
4119   }
4120   }
4121   }
4122   }
4123   }
4124   }
4125   }
4126   }
4127   }
4128   }
4129   }
4130   }
4131   }
4132   }
4133   }
4134   }
4135   }
4136   }
4137   }
4138   }
4139   }
4140   }
4141   }
4142   }
4143   }
4144   }
4145   }
4146   }
4147   }
4148   }
4149   }
4150   }
4151   }
4152   }
4153   }
4154   }
4155   }
4156   }
4157   }
4158   }
4159   }
4160   }
4161   }
4162   }
4163   }
4164   }
4165   }
4166   }
4167   }
4168   }
4169   }
4170   }
4171   }
4172   }
4173   }
4174   }
4175   }
4176   }
4177   }
4178   }
4179   }
4180   }
4181   }
4182   }
4183   }
4184   }
4185   }
4186   }

```



```

4376 * fpropen(oddnum, dskm, w);
4377 * ftime(tp); % of MEDIA.MAN clock, tr).
4378 */
4379 tr_phase [ tr ] = PRZ;
4380 cc_free(tr);
4381 cc_free(tr);
4382
4383 if (class_of [ tr ] == REMOTE) {
4384   broadcast ( NEED_BUS_TOKEN, PRE_TIME, token );
4385   (void) ftime(tp);
4386   free [ tr ];
4387   free [ tr ];
4388   num_alive [ tr ] = 0;
4389   num_in_q [ tr ] = 0;
4390   blocked_site [ tr ] = 0;
4391   free [ tr ];
4392   free [ tr ];
4393   free [ tr ];
4394   free [ tr ];
4395   free [ tr ];
4396   goto state ( LP, PRE_LOCK );
4397   goto phase ( LP );
4398   setUp ( token );
4399 }
4400
4401 /****** initialize site relation card *****/
4402 for_all_rel for_all_site
4403   if ( queried_rel [ tr ] [ rel ] == 1 )
4404     rel_card [ rel ] [ site ] [ rel ] = 1;
4405   rel_card [ rel ] [ site ] [ rel ] = 1;
4406   schedule ( NEED_TR, spact (1000000 * mean_tr), 0 );
4407 }
4408 }
4409 }
4410 }
4411 }
4412 /* calculate effective relation width ( sum of queried site )
4413 for_all_rel if ( queried_rel [ tr ] [ rel ] == 1 ) {
4414   tr_rel_width [ tr ] [ rel ] = 0;
4415   for_all_queried_site_in_rel
4416     for_all_queried_site_in_rel
4417       if ( queried_rel [ tr ] [ rel ] == 1 ) {
4418         tr_rel_width [ tr ] [ rel ] += relation_card [ rel ];
4419         tr_rel_width [ tr ] [ rel ] += relation_card [ rel ];
4420       }
4421 }
4422 /* set up effective relation width for each relation */
4423 for_all_rel if ( queried_rel [ tr ] [ rel ] == 1 ) {
4424   tr_rel_width [ tr ] [ rel ] = 0;
4425   for_all_queried_site_in_rel
4426     for_all_queried_site_in_rel
4427       if ( queried_rel [ tr ] [ rel ] == 1 ) {
4428         tr_rel_width [ tr ] [ rel ] += relation_card [ rel ];
4429         tr_rel_width [ tr ] [ rel ] += relation_card [ rel ];
4430       }
4431 }
4432 }
4433 }
4434 }
4435 }
4436 }
4437 }
4438 }
4439 }
4440 }
4441 * if home site has some queried relations, the condition is set to false
4442 * since test join process is needed in this case
4443 * when data is an orphan site, the site will perform join in the site
4444 * then send the results directly, no phase is desired
4445 */
4446 orphan_test ( token )
4447 {
4448   test;
4449 }
4450 }
4451 }
4452 }
4453 }
4454 }
4455 }
4456 }
4457 }
4458 }
4459 }
4460 }
4461 }
4462 }
4463 }
4464 }
4465 }
4466 }
4467 }
4468 }
4469 }
4470 }
4471 }
4472 }
4473 }
4474 }
4475 }
4476 }
4477 }
4478 }
4479 }
4480 }
4481 }
4482 }
4483 }
4484 }
4485 }
4486 }
4487 }
4488 }
4489 }
4490 }
4491 }
4492 }
4493 }
4494 }
4495 }
4496 }
4497 }
4498 }
4499 }
4500 }
4501 }
4502 }
4503 }
4504 }
4505 }
4506 }
4507 }
4508 }
4509 }
4510 }
4511 }
4512 }
4513 }
4514 }
4515 }
4516 }
4517 }
4518 }
4519 }
4520 }
4521 }
4522 }
4523 }
4524 }
4525 }
4526 }
4527 }
4528 }
4529 }
4530 }
4531 }
4532 }
4533 }
4534 }
4535 }
4536 }
4537 }
4538 }
4539 }
4540 }
4541 }
4542 }
4543 }
4544 }
4545 }
4546 }
4547 }
4548 }
4549 }
4550 }
4551 }
4552 }
4553 }
4554 }
4555 }
4556 }
4557 }
4558 }
4559 }
4560 }
4561 }
4562 }
4563 }
4564 }
4565 }
4566 }
4567 }
4568 }
4569 }
4570 }
4571 }
4572 }
4573 }
4574 }
4575 }
4576 }
4577 }
4578 }
4579 }
4580 }
4581 }
4582 }
4583 }
4584 }
4585 }
4586 }
4587 }
4588 }
4589 }
4590 }
4591 }
4592 }
4593 }
4594 }
4595 }
4596 }
4597 }
4598 }
4599 }
4600 }
4601 }
4602 }
4603 }
4604 }
4605 }
4606 }
4607 }
4608 }
4609 }
4610 }
4611 }
4612 }
4613 }
4614 }
4615 }
4616 }
4617 }
4618 }
4619 }
4620 }
4621 }
4622 }
4623 }
4624 }
4625 }
4626 }
4627 }
4628 }
4629 }
4630 }
4631 }
4632 }
4633 }
4634 }
4635 }
4636 }
4637 }
4638 }
4639 }
4640 }
4641 }
4642 }
4643 }
4644 }
4645 }
4646 }
4647 }
4648 }
4649 }
4650 }
4651 }
4652 }
4653 }
4654 }
4655 }
4656 }
4657 }
4658 }
4659 }
4660 }
4661 }
4662 }
4663 }
4664 }
4665 }
4666 }
4667 }
4668 }
4669 }
4670 }
4671 }
4672 }
4673 }
4674 }
4675 }
4676 }
4677 }
4678 }
4679 }
4680 }
4681 }
4682 }
4683 }
4684 }
4685 }
4686 }
4687 }
4688 }
4689 }
4690 }
4691 }
4692 }
4693 }
4694 }
4695 }
4696 }
4697 }
4698 }
4699 }
4700 }
4701 }
4702 }
4703 }
4704 }
4705 }
4706 }
4707 }
4708 }
4709 }
4710 }
4711 }
4712 }
4713 }
4714 }
4715 }
4716 }
4717 }
4718 }
4719 }
4720 }
4721 }
4722 }
4723 }
4724 }
4725 }
4726 }
4727 }
4728 }
4729 }
4730 }
4731 }
4732 }
4733 }
4734 }
4735 }
4736 }
4737 }
4738 }
4739 }
4740 }
4741 }
4742 }
4743 }
4744 }
4745 }
4746 }
4747 }
4748 }
4749 }
4750 }
4751 }
4752 }
4753 }
4754 }
4755 }
4756 }
4757 }
4758 }
4759 }
4760 }
4761 }
4762 }
4763 }
4764 }
4765 }
4766 }
4767 }
4768 }
4769 }
4770 }
4771 }
4772 }
4773 }
4774 }
4775 }
4776 }
4777 }
4778 }
4779 }
4780 }
4781 }
4782 }
4783 }
4784 }
4785 }
4786 }
4787 }
4788 }
4789 }
4790 }
4791 }
4792 }
4793 }
4794 }
4795 }
4796 }
4797 }
4798 }
4799 }
4800 }
4801 }
4802 }
4803 }
4804 }
4805 }
4806 }
4807 }
4808 }
4809 }
4810 }
4811 }
4812 }
4813 }
4814 }
4815 }
4816 }
4817 }
4818 }
4819 }
4820 }
4821 }
4822 }
4823 }
4824 }
4825 }
4826 }
4827 }
4828 }
4829 }
4830 }
4831 }
4832 }
4833 }
4834 }
4835 }
4836 }
4837 }
4838 }
4839 }
4840 }
4841 }
4842 }
4843 }
4844 }
4845 }
4846 }
4847 }
4848 }
4849 }
4850 }
4851 }
4852 }
4853 }
4854 }
4855 }
4856 }
4857 }
4858 }
4859 }
4860 }
4861 }
4862 }
4863 }
4864 }
4865 }
4866 }
4867 }
4868 }
4869 }
4870 }
4871 }
4872 }
4873 }
4874 }
4875 }
4876 }
4877 }
4878 }
4879 }
4880 }
4881 }
4882 }
4883 }
4884 }
4885 }
4886 }
4887 }
4888 }
4889 }
4890 }
4891 }
4892 }
4893 }
4894 }
4895 }
4896 }
4897 }
4898 }
4899 }
4900 }
4901 }
4902 }
4903 }
4904 }
4905 }
4906 }
4907 }
4908 }
4909 }
4910 }
4911 }
4912 }
4913 }
4914 }
4915 }
4916 }
4917 }
4918 }
4919 }
4920 }
4921 }
4922 }
4923 }
4924 }
4925 }
4926 }
4927 }
4928 }
4929 }
4930 }
4931 }
4932 }
4933 }
4934 }
4935 }
4936 }
4937 }
4938 }
4939 }
4940 }
4941 }
4942 }
4943 }
4944 }
4945 }
4946 }
4947 }
4948 }
4949 }
4950 }
4951 }
4952 }
4953 }
4954 }
4955 }
4956 }
4957 }
4958 }
4959 }
4960 }
4961 }
4962 }
4963 }
4964 }
4965 }
4966 }
4967 }
4968 }
4969 }
4970 }
4971 }
4972 }
4973 }
4974 }
4975 }
4976 }
4977 }
4978 }
4979 }
4980 }
4981 }
4982 }
4983 }
4984 }
4985 }
4986 }
4987 }
4988 }
4989 }
4990 }
4991 }
4992 }
4993 }
4994 }
4995 }
4996 }
4997 }
4998 }
4999 }
5000 }

```

```

4460 }
4461 }
4462 }
4463 }
4464 }
4465 }
4466 }
4467 }
4468 }
4469 }
4470 }
4471 }
4472 }
4473 }
4474 }
4475 }
4476 }
4477 }
4478 }
4479 }
4480 }
4481 }
4482 }
4483 }
4484 }
4485 }
4486 }
4487 }
4488 }
4489 }
4490 }
4491 }
4492 }
4493 }
4494 }
4495 }
4496 }
4497 }
4498 }
4499 }
4500 }
4501 }
4502 }
4503 }
4504 }
4505 }
4506 }
4507 }
4508 }
4509 }
4510 }
4511 }
4512 }
4513 }
4514 }
4515 }
4516 }
4517 }
4518 }
4519 }
4520 }
4521 }
4522 }
4523 }
4524 }
4525 }
4526 }
4527 }
4528 }
4529 }
4530 }
4531 }
4532 }
4533 }
4534 }
4535 }
4536 }
4537 }
4538 }
4539 }
4540 }
4541 }
4542 }
4543 }
4544 }
4545 }
4546 }
4547 }
4548 }
4549 }
4550 }
4551 }
4552 }
4553 }
4554 }
4555 }
4556 }
4557 }
4558 }
4559 }
4560 }
4561 }
4562 }
4563 }
4564 }
4565 }
4566 }
4567 }
4568 }
4569 }
4570 }
4571 }
4572 }
4573 }
4574 }
4575 }
4576 }
4577 }
4578 }
4579 }
4580 }
4581 }
4582 }
4583 }
4584 }
4585 }
4586 }
4587 }
4588 }
4589 }
4590 }
4591 }
4592 }
4593 }
4594 }
4595 }
4596 }
4597 }
4598 }
4599 }
4600 }
4601 }
4602 }
4603 }
4604 }
4605 }
4606 }
4607 }
4608 }
4609 }
4610 }
4611 }
4612 }
4613 }
4614 }
4615 }
4616 }
4617 }
4618 }
4619 }
4620 }
4621 }
4622 }
4623 }
4624 }
4625 }
4626 }
4627 }
4628 }
4629 }
4630 }
4631 }
4632 }
4633 }
4634 }
4635 }
4636 }
4637 }
4638 }
4639 }
4640 }
4641 }
4642 }
4643 }
4644 }
4645 }
4646 }
4647 }
4648 }
4649 }
4650 }
4651 }
4652 }
4653 }
4654 }
4655 }
4656 }
4657 }
4658 }
4659 }
4660 }
4661 }
4662 }
4663 }
4664 }
4665 }
4666 }
4667 }
4668 }
4669 }
4670 }
4671 }
4672 }
4673 }
4674 }
4675 }
4676 }
4677 }
4678 }
4679 }
4680 }
4681 }
4682 }
4683 }
4684 }
4685 }
4686 }
4687 }
4688 }
4689 }
4690 }
4691 }
4692 }
4693 }
4694 }
4695 }
4696 }
4697 }
4698 }
4699 }
4700 }
4701 }
4702 }
4703 }
4704 }
4705 }
4706 }
4707 }
4708 }
4709 }
4710 }
4711 }
4712 }
4713 }
4714 }
4715 }
4716 }
4717 }
4718 }
4719 }
4720 }
4721 }
4722 }
4723 }
4724 }
4725 }
4726 }
4727 }
4728 }
4729 }
4730 }
4731 }
4732 }
4733 }
4734 }
4735 }
4736 }
4737 }
4738 }
4739 }
4740 }
4741 }
4742 }
4743 }
4744 }
4745 }
4746 }
4747 }
4748 }
4749 }
4750 }
4751 }
4752 }
4753 }
4754 }
4755 }
4756 }
4757 }
4758 }
4759 }
4760 }
4761 }
4762 }
4763 }
4764 }
4765 }
4766 }
4767 }
4768 }
4769 }
4770 }
4771 }
4772 }
4773 }
4774 }
4775 }
4776 }
4777 }
4778 }
4779 }
4780 }
4781 }
4782 }
4783 }
4784 }
4785 }
4786 }
4787 }
4788 }
4789 }
4790 }
4791 }
4792 }
4793 }
4794 }
4795 }
4796 }
4797 }
4798 }
4799 }
4800 }
4801 }
4802 }
4803 }
4804 }
4805 }
4806 }
4807 }
4808 }
4809 }
4810 }
4811 }
4812 }
4813 }
4814 }
4815 }
4816 }
4817 }
4818 }
4819 }
4820 }
4821 }
4822 }
4823 }
4824 }
4825 }
4826 }
4827 }
4828 }
4829 }
4830 }
4831 }
4832 }
4833 }
4834 }
4835 }
4836 }
4837 }
4838 }
4839 }
4840 }
4841 }
4842 }
4843 }
4844 }
4845 }
4846 }
4847 }
4848 }
4849 }
4850 }
4851 }
4852 }
4853 }
4854 }
4855 }
4856 }
4857 }
4858 }
4859 }
4860 }
4861 }
4862 }
4863 }
4864 }
4865 }
4866 }
4867 }
4868 }
4869 }
4870 }
4871 }
4872 }
4873 }
4874 }
4875 }
4876 }
4877 }
4878 }
4879 }
4880 }
4881 }
4882 }
4883 }
4884 }
4885 }
4886 }
4887 }
4888 }
4889 }
4890 }
4891 }
4892 }
4893 }
4894 }
4895 }
4896 }
4897 }
4898 }
4899 }
4900 }
4901 }
4902 }
4903 }
4904 }
4905 }
4906 }
4907 }
4908 }
4909 }
4910 }
4911 }
4912 }
4913 }
4914 }
4915 }
4916 }
4917 }
4918 }
4919 }
4920 }
4921 }
4922 }
4923 }
4924 }
4925 }
4926 }
4927 }
4928 }
4929 }
4930 }
4931 }
4932 }
4933 }
4934 }
4935 }
4936 }
4937 }
4938 }
4939 }
4940 }
4941 }
4942 }
4943 }
4944 }
4945 }
4946 }
4947 }
4948 }
4949 }
4950 }
4951 }
4952 }
4953 }
4954 }
4955 }
4956 }
4957 }
4958 }
4959 }
4960 }
4961 }
4962 }
4963 }
4964 }
4965 }
4966 }
4967 }
4968 }
4969 }
4970 }
4971 }
4972 }
4973 }
4974 }
4975 }
4976 }
4977 }
4978 }
4979 }
4980 }
4981 }
4982 }
4983 }
4984 }
4985 }
4986 }
4987 }
4988 }
4989 }
4990 }
4991 }
4992 }
4993 }
4994 }
4995 }
4996 }
4997 }
4998 }
4999 }
5000 }

```

```

4503         alive [ i ] = 0 ;
4504         goto loop ( tr, i, OBJ_WAIT_GJ ) ;
4505     }
4506     }
4507     /* Had picked up a the min. site */
4508     /* Had picked up a the min. site */
4509     /* Had picked up a the min. site */
4510     /* Had picked up a the min. site */
4511     /* Had picked up a the min. site */
4512     /* Had picked up a the min. site */
4513     /* Had picked up a the min. site */
4514     /* Had picked up a the min. site */
4515     /* Had picked up a the min. site */
4516     /* Had picked up a the min. site */
4517     /* Had picked up a the min. site */
4518     /* Had picked up a the min. site */
4519     /* Had picked up a the min. site */
4520     /* Had picked up a the min. site */
4521     /* Had picked up a the min. site */
4522     /* Had picked up a the min. site */
4523     /* Had picked up a the min. site */
4524     /* Had picked up a the min. site */
4525     /* Had picked up a the min. site */
4526     /* Had picked up a the min. site */
4527     /* Had picked up a the min. site */
4528     /* Had picked up a the min. site */
4529     /* Had picked up a the min. site */
4530     /* Had picked up a the min. site */
4531     /* Had picked up a the min. site */
4532     /* Had picked up a the min. site */
4533     /* Had picked up a the min. site */
4534     /* Had picked up a the min. site */
4535     /* Had picked up a the min. site */
4536     /* Had picked up a the min. site */
4537     /* Had picked up a the min. site */
4538     /* Had picked up a the min. site */
4539     /* Had picked up a the min. site */
4540     /* Had picked up a the min. site */
4541     /* Had picked up a the min. site */
4542     /* Had picked up a the min. site */
4543     /* Had picked up a the min. site */
4544     /* Had picked up a the min. site */
4545     /* Had picked up a the min. site */
4546     /* Had picked up a the min. site */
4547     /* Had picked up a the min. site */
4548     /* Had picked up a the min. site */
4549     /* Had picked up a the min. site */
4550     /* Had picked up a the min. site */
4551     /* Had picked up a the min. site */
4552     /* Had picked up a the min. site */
4553     /* Had picked up a the min. site */
4554     /* Had picked up a the min. site */
4555     /* Had picked up a the min. site */
4556     /* Had picked up a the min. site */
4557     /* Had picked up a the min. site */
4558     /* Had picked up a the min. site */
4559     /* Had picked up a the min. site */
4560     /* Had picked up a the min. site */
4561     /* Had picked up a the min. site */
4562     /* Had picked up a the min. site */
4563     /* Had picked up a the min. site */
4564     /* Had picked up a the min. site */
4565     /* Had picked up a the min. site */
4566     /* Had picked up a the min. site */

```

```

4567     queried_site [tr] [rel] [act] = 1.
4568     }
4569     }
4570     }
4571     }
4572     }
4573     }
4574     }
4575     }
4576     }
4577     }
4578     }
4579     }
4580     }
4581     }
4582     }
4583     }
4584     }
4585     }
4586     }
4587     }
4588     }
4589     }
4590     }
4591     }
4592     }
4593     }
4594     }
4595     }
4596     }
4597     }
4598     }
4599     }
4600     }
4601     }
4602     }
4603     }
4604     }
4605     }
4606     }
4607     }
4608     }
4609     }
4610     }
4611     }
4612     }
4613     }
4614     }
4615     }
4616     }
4617     }
4618     }
4619     }
4620     }
4621     }
4622     }
4623     }
4624     }
4625     }
4626     }
4627     }
4628     }
4629     }
4630     }
4631     }
4632     }
4633     }
4634     }
4635     }
4636     }
4637     }
4638     }
4639     }
4640     }
4641     }
4642     }
4643     }
4644     }
4645     }
4646     }
4647     }
4648     }
4649     }
4650     }
4651     }
4652     }
4653     }
4654     }
4655     }
4656     }
4657     }
4658     }
4659     }
4660     }
4661     }
4662     }
4663     }
4664     }
4665     }
4666     }
4667     }
4668     }
4669     }
4670     }
4671     }
4672     }
4673     }
4674     }
4675     }
4676     }
4677     }
4678     }
4679     }
4680     }
4681     }
4682     }
4683     }
4684     }
4685     }
4686     }
4687     }
4688     }
4689     }
4690     }
4691     }
4692     }
4693     }
4694     }
4695     }
4696     }
4697     }
4698     }
4699     }
4700     }
4701     }
4702     }
4703     }
4704     }
4705     }
4706     }
4707     }
4708     }
4709     }
4710     }
4711     }
4712     }
4713     }
4714     }
4715     }
4716     }
4717     }
4718     }
4719     }
4720     }
4721     }
4722     }
4723     }
4724     }
4725     }
4726     }
4727     }
4728     }
4729     }
4730     }
4731     }
4732     }
4733     }
4734     }
4735     }
4736     }
4737     }
4738     }
4739     }
4740     }
4741     }
4742     }
4743     }
4744     }
4745     }
4746     }
4747     }
4748     }
4749     }
4750     }
4751     }
4752     }
4753     }
4754     }
4755     }
4756     }
4757     }
4758     }
4759     }
4760     }

```

```

4761 */
4762 if ( queried_frel [tr] [rel] == 1 ) {
4763     i=0;
4764     for_all_act
4765     if ( queried_act [tr] [rel] [act] == 1 ) { i++; break; }
4766     if (i==0) {
4767         j=RANDOM(R);
4768         while ( j > 0 ) {
4769             act = ( (k-1) * J.A ) * i;
4770             if ( relation_conf [rel] [act] == 1 ) j--;
4771             K++;
4772         }
4773     }
4774     queried_act [tr] [rel] [act] = i;
4775 }
4776 }
4777 /* Test for the existence of joining act
4778 */
4779 /* Either zero or more than one should be stated
4780 */
4781 for_all_act {
4782     i=0;
4783     for_all_frel if ( queried_act [tr] [rel] [act] == 1 ) i++;
4784     if ( i == 1 ) {
4785         i=0;
4786         while ( i == 0 ) {
4787             i = RANDOM(R);
4788             relation_conf [i] [act] == 0) i=0;
4789         }
4790         queried_act [tr] [i] [act] = queried_frel [tr] [i] = 1;
4791     }
4792 }
4793 /*
4794 * resolve_local_id () --- Resolving possible local deadlock
4795 *
4796 * CALLED BY LP_WAIT_SVN on finish
4797 *
4798 * If the estimated slowest site is blocked by somebody else, the
4799 * local pg held by the TR may cause deadlock problems
4800 * Changing the est. slowest site selection to prevent it
4801 *
4802 * Check if the estimated slowest site is blocked. If YES then do it
4803 * Assign the est. slowest site the one of the sites in state LP_DISK
4804 * If no such site, then choose a site in LP_WAIT_END. Make it up
4805 */
4806 est_slowest_site [tr] should be changed
4807
4808 4809
4809     reactive_local_id ( token )
4810     for token;
4811     register int KK, tr;
4812     int i, j, rel, site, act;
4813     site = site_of [ token ];
4814     tr = tr_of [ token ];
4815     pre ( resolve_local_id );
4816
4817     KK = est_slowest_site_of [ tr ];
4818     if (
4819         KK > 0
4820         || KK < INF
4821         || LP_WAIT_SVN
4822     ) {
4823         /*
4824         * SVN had been broadcast, do local alone
4825         * do not broadcast token. (reactonly)
4826         */
4827         return;
4828     }
4829     if ( cc_lock_get_busy ( tr, site ) == FAIL ) {
4830         return ( FAIL );
4831     }
4832     return ( FAIL );
4833     reactive_local_id and cc_lock_get_busy );
4834 }
4835 }
4836 }
4837 }
4838 }
4839 }
4840 }
4841 }
4842 }
4843 }
4844 }
4845 }
4846 }
4847 }
4848 }
4849 }
4850 }
4851 }
4852 }
4853 }
4854 }
4855 }
4856 }
4857 }
4858 }
4859 }
4860 }
4861 }
4862 }
4863 }
4864 }
4865 }
4866 }
4867 }
4868 }
4869 }
4870 }
4871 }
4872 }
4873 }
4874 }
4875 }
4876 }
4877 }
4878 }
4879 }
4880 }
4881 }
4882 }
4883 }
4884 }
4885 }
4886 }
4887 }
4888 }
4889 }
4890 }
4891 }
4892 }
4893 }
4894 }
4895 }
4896 }
4897 }
4898 }
4899 }
4900 }
4901 }
4902 }
4903 }
4904 }
4905 }
4906 }
4907 }
4908 }
4909 }
4910 }
4911 }
4912 }
4913 }
4914 }
4915 }
4916 }
4917 }
4918 }
4919 }
4920 }
4921 }
4922 }
4923 }
4924 }
4925 }
4926 }
4927 }
4928 }
4929 }
4930 }
4931 }
4932 }
4933 }
4934 }
4935 }
4936 }
4937 }
4938 */

```



```

6304 *
6305 * unlock
6306 * changed BPT
6307 * (change blocked processes)
6308 * delete the transaction, return resources back
6309 */
6310
6311 terminate_site ( token, tr, site )
6312 int token;
6313 register int tr, site;
6314
6315 /* Pre ( terminate_site ) */
6316
6317 /* Post ( terminate_site ) */
6318
6319
6320
6321
6322
6323
6324
6325
6326
6327
6328
6329
6330
6331
6332
6333
6334
6335
6336
6337
6338
6339
6340
6341
6342
6343
6344
6345
6346
6347
6348
6349
6350
6351
6352
6353
6354
6355
6356
6357
6358
6359
6360
6361
6362
6363
6364
6365
6366
6367
6368
6369
6370
6371
6372
6373
6374
6375
6376
6377
6378
6379
6380
6381
6382
6383
6384
6385
6386
6387
6388
6389
6390
6391
6392
6393
6394
6395
6396
6397

```

```

6398
6399
6400
6401
6402
6403
6404
6405
6406
6407
6408
6409
6410
6411
6412
6413
6414
6415
6416
6417
6418
6419
6420
6421
6422
6423
6424
6425
6426
6427
6428
6429
6430
6431
6432
6433
6434
6435
6436
6437
6438
6439
6440
6441
6442
6443
6444
6445
6446
6447
6448
6449
6450
6451
6452
6453
6454
6455
6456
6457
6458
6459
6460
6461
6462
6463
6464
6465
6466
6467
6468
6469
6470
6471
6472
6473
6474
6475
6476
6477
6478
6479
6480
6481
6482
6483
6484
6485
6486
6487
6488
6489
6490
6491
6492
6493
6494
6495
6496
6497
6498
6499
6500
6501

```


6233	Report	1196 6708 6708 6722 6836 6847
6234	Reporter	6708 6722 6836 6847
6236	Reserved	804 1129 1136 1167 6002 6010 6012
6238	Reset	5716
6240	Reset	4761
6241	Resolve_Local_d1	
6242		
6243	Resume_By	1638 1663 1672 1621 1736 1807 1819 2319 2389 4833 4871 4876
6246	Return_Locked_Id	1609 4913 6316
6247	Return_Tr_Id	4933 6388
6248	Tr_Conflict	3208 3287 4964
6250		
6261	Schedule	481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000
6262	Size_of_Fl_Bro	1281 1366 6034 6040
6263	Smpl	683 722 724 726 736 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000
6264	Stat	1047 6676
6265	Stat	1002 6034
6266	Stat	1047 6676
6267	Stat	1002 6034
6268	Stat	1047 6676
6269	Stat	1002 6034
6270	Stat	1047 6676
6271	Stat	1002 6034
6272	Stat	1047 6676
6273	Stat	1002 6034
6274	Stat	1047 6676
6275	Stat	1002 6034
6276	Stat	1047 6676
6277	Stat	1002 6034
6278	Stat	1047 6676
6279	Stat	1002 6034
6280	Stat	1047 6676
6281	Stat	1002 6034
6282	Stat	1047 6676
6283	Stat	1002 6034
6284	Stat	1047 6676
6285	Stat	1002 6034
6286	Stat	1047 6676
6287	Stat	1002 6034
6288	Stat	1047 6676
6289	Stat	1002 6034
6290	Stat	1047 6676
6291	Stat	1002 6034
6292	Stat	1047 6676
6293	Stat	1002 6034
6294	Stat	1047 6676
6295	Stat	1002 6034
6296	Stat	1047 6676
6297	Stat	1002 6034
6298	Stat	1047 6676
6299	Stat	1002 6034
6300	Stat	1047 6676
6301	Stat	1002 6034
6302	Stat	1047 6676
6303	Stat	1002 6034
6304	Stat	1047 6676
6305	Stat	1002 6034
6306	Stat	1047 6676
6307	Stat	1002 6034
6308	Stat	1047 6676
6309	Stat	1002 6034
6310	Stat	1047 6676
6311	Stat	1002 6034
6312	Stat	1047 6676
6313	Stat	1002 6034
6314	Stat	1047 6676
6315	Stat	1002 6034
6316	Stat	1047 6676
6317	Stat	1002 6034
6318	Stat	1047 6676
6319	Stat	1002 6034
6320	Stat	1047 6676
6321	Stat	1002 6034
6322	Stat	1047 6676
6323	Stat	1002 6034
6324	Stat	1047 6676
6325	Stat	1002 6034
6326	Stat	1047 6676
6327	Stat	1002 6034
6328	Stat	1047 6676
6329	Stat	1002 6034
6330	Stat	1047 6676
6331	Stat	1002 6034
6332	Stat	1047 6676
6333	Stat	1002 6034
6334	Stat	1047 6676
6335	Stat	1002 6034
6336	Stat	1047 6676
6337	Stat	1002 6034
6338	Stat	1047 6676
6339	Stat	1002 6034
6340	Stat	1047 6676
6341	Stat	1002 6034
6342	Stat	1047 6676
6343	Stat	1002 6034
6344	Stat	1047 6676
6345	Stat	1002 6034
6346	Stat	1047 6676
6347	Stat	1002 6034
6348	Stat	1047 6676
6349	Stat	1002 6034
6350	Stat	1047 6676
6351	Stat	1002 6034
6352	Stat	1047 6676
6353	Stat	1002 6034
6354	Stat	1047 6676
6355	Stat	1002 6034
6356	Stat	1047 6676
6357	Stat	1002 6034
6358	Stat	1047 6676
6359	Stat	1002 6034
6360	Stat	1047 6676
6361	Stat	1002 6034
6362	Stat	1047 6676
6363	Stat	1002 6034
6364	Stat	1047 6676
6365	Stat	1002 6034
6366	Stat	1047 6676
6367	Stat	1002 6034
6368	Stat	1047 6676
6369	Stat	1002 6034
6370	Stat	1047 6676
6371	Stat	1002 6034
6372	Stat	1047 6676
6373	Stat	1002 6034
6374	Stat	1047 6676
6375	Stat	1002 6034
6376	Stat	1047 6676
6377	Stat	1002 6034
6378	Stat	1047 6676
6379	Stat	1002 6034
6380	Stat	1047 6676
6381	Stat	1002 6034
6382	Stat	1047 6676
6383	Stat	1002 6034
6384	Stat	1047 6676
6385	Stat	1002 6034
6386	Stat	1047 6676
6387	Stat	1002 6034
6388	Stat	1047 6676
6389	Stat	1002 6034
6390	Stat	1047 6676
6391	Stat	1002 6034
6392	Stat	1047 6676
6393	Stat	1002 6034
6394	Stat	1047 6676
6395	Stat	1002 6034
6396	Stat	1047 6676
6397	Stat	1002 6034
6398	Stat	1047 6676
6399	Stat	1002 6034
6400	Stat	1047 6676
6401	Stat	1002 6034
6402	Stat	1047 6676
6403	Stat	1002 6034
6404	Stat	1047 6676
6405	Stat	1002 6034
6406	Stat	1047 6676
6407	Stat	1002 6034
6408	Stat	1047 6676
6409	Stat	1002 6034
6410	Stat	1047 6676
6411	Stat	1002 6034
6412	Stat	1047 6676
6413	Stat	1002 6034
6414	Stat	1047 6676
6415	Stat	1002 6034
6416	Stat	1047 6676
6417	Stat	1002 6034
6418	Stat	1047 6676
6419	Stat	1002 6034
6420	Stat	1047 6676
6421	Stat	1002 6034
6422	Stat	1047 6676
6423	Stat	1002 6034
6424	Stat	1047 6676
6425	Stat	1002 6034
6426	Stat	1047 6676
6427	Stat	1002 6034
6428	Stat	1047 6676
6429	Stat	1002 6034
6430	Stat	1047 6676

Cross Reference of Important Global Variables and Constants

DATA SIZE

1099 6340 6342 316 1099 1347 1428 6064

0421 char 619 627 636 694 706 708 5848 8661 8741 8790 8821 8864 8961 8970 8974
0422 stringbuf 708 6964 6967 6979
0423 class_of 343 1368 1822 2001 2006 2276 2349 2442 3072 3704 3767 3774 3778 3790
0424 component_map 344 608 1032 1034 3610 3922 4042 8413
0425 counter 284 481 494 496 499
0426 disk 303 322 346 412 466 1088 1167 1168 1169 1184 1848 2038 2039 2047 3947
0427 domain_range 3980 4368 6346 6381 6883
0428 dir 322 2038 2039
0429 domain_width 696 1341 3460 3482 3902 4007 4210
0430 end_line 680 1342 1422 1909 1938 1944 3902 4022 4208 4410 4420 8049
0431 end_line 286 1080 1192 1193
0432 end_line 286 1080 1192 1193
0433 end_line 346 3440 3481 4207 4208
0434 end_line 347 3436 3448 3494 3495 3939
0435 end_line_of 348 2081 2094 2459 2656 4391 4861 4866 6080 6143 6366
0436 event 272 683 741 747 967 968 1081 1109 1113 1118 1121 6073 6123 6126
0437 association_line 6126 6128 6616 6617 6650 6746 8017
0438 association_line 287 2463 2464
0439 tap 466 476 479 489 491 493 609 1028 1096 1412 1704 2023 2937 3470
0440 streambuf 3468 3627 3850 3788 3796 3807 4289 4289
0441 streambuf 276 1108 3946 6368 6370 6373 6881 8862 8883 8884
0442 allowed_error 349 2630 2633 3871 4380 4386 6384 6389
0443 allowed_error 306 1070 1106 6374
0444 allowed_error 360 1726 1786 2620 3110 3828 3768 3777 3796 4079 4366 4943 6346 6383
0445 allowed_error 361 361 1883 3132 3133 3878 3910 4196 4200 4207 4212 4213 4235 4239
0446 allowed_error 4240 4241 4628 4662 6321
0447 allowed_error 362 684 1643 1976 3162 3460 4000 4001 4007 4013 4399
0448 allowed_error 684 4001
0449 iteration 328 1080 1089 1106 6371 6374
0450 iteration 328 1108 6371 6374
0451 iteration 328 1080 1089 1106 6371 6374
0452 iteration 328 1108 6371 6374
0453 iteration 328 1080 1089 1106 6371 6374
0454 iteration 328 1080 1089 1106 6371 6374
0455 iteration 328 1080 1089 1106 6371 6374
0456 iteration 328 1080 1089 1106 6371 6374
0457 iteration 328 1080 1089 1106 6371 6374
0458 iteration 328 1080 1089 1106 6371 6374
0459 iteration 328 1080 1089 1106 6371 6374
0460 iteration 328 1080 1089 1106 6371 6374
0461 iteration 328 1080 1089 1106 6371 6374
0462 iteration 328 1080 1089 1106 6371 6374
0463 iteration 328 1080 1089 1106 6371 6374
0464 iteration 328 1080 1089 1106 6371 6374
0465 iteration 328 1080 1089 1106 6371 6374
0466 iteration 328 1080 1089 1106 6371 6374
0467 iteration 328 1080 1089 1106 6371 6374
0468 iteration 328 1080 1089 1106 6371 6374
0469 iteration 328 1080 1089 1106 6371 6374
0470 iteration 328 1080 1089 1106 6371 6374
0471 iteration 328 1080 1089 1106 6371 6374
0472 iteration 328 1080 1089 1106 6371 6374
0473 iteration 328 1080 1089 1106 6371 6374
0474 iteration 328 1080 1089 1106 6371 6374
0475 iteration 328 1080 1089 1106 6371 6374
0476 iteration 328 1080 1089 1106 6371 6374
0477 iteration 328 1080 1089 1106 6371 6374
0478 iteration 328 1080 1089 1106 6371 6374
0479 iteration 328 1080 1089 1106 6371 6374
0480 iteration 328 1080 1089 1106 6371 6374
0481 iteration 328 1080 1089 1106 6371 6374
0482 iteration 328 1080 1089 1106 6371 6374
0483 iteration 328 1080 1089 1106 6371 6374
0484 iteration 328 1080 1089 1106 6371 6374
0485 iteration 328 1080 1089 1106 6371 6374
0486 iteration 328 1080 1089 1106 6371 6374
0487 iteration 328 1080 1089 1106 6371 6374
0488 iteration 328 1080 1089 1106 6371 6374
0489 iteration 328 1080 1089 1106 6371 6374
0490 iteration 328 1080 1089 1106 6371 6374
0491 iteration 328 1080 1089 1106 6371 6374
0492 iteration 328 1080 1089 1106 6371 6374
0493 iteration 328 1080 1089 1106 6371 6374
0494 iteration 328 1080 1089 1106 6371 6374
0495 iteration 328 1080 1089 1106 6371 6374
0496 iteration 328 1080 1089 1106 6371 6374
0497 iteration 328 1080 1089 1106 6371 6374
0498 iteration 328 1080 1089 1106 6371 6374
0499 iteration 328 1080 1089 1106 6371 6374
0500 iteration 328 1080 1089 1106 6371 6374
0501 iteration 328 1080 1089 1106 6371 6374
0502 iteration 328 1080 1089 1106 6371 6374
0503 iteration 328 1080 1089 1106 6371 6374
0504 iteration 328 1080 1089 1106 6371 6374
0505 iteration 328 1080 1089 1106 6371 6374
0506 iteration 328 1080 1089 1106 6371 6374
0507 iteration 328 1080 1089 1106 6371 6374
0508 iteration 328 1080 1089 1106 6371 6374
0509 iteration 328 1080 1089 1106 6371 6374
0510 iteration 328 1080 1089 1106 6371 6374
0511 iteration 328 1080 1089 1106 6371 6374
0512 iteration 328 1080 1089 1106 6371 6374
0513 iteration 328 1080 1089 1106 6371 6374
0514 iteration 328 1080 1089 1106 6371 6374

0515 new_inactive 384 2287 2366 2790 2724 2736 2738 2742 2749 2860 2863 2865 2866
0516 new_inactive 3198 3334 3248 3258 3267 3240 3241 3251 3262 3263 3264 3265 3266
0517 new_inactive 3371 3272 3276 3286 3324 3331 3332 3337 3368 3387 3379 3380
0518 new_inactive 4887 4888 6100 6141 6348 6419
0519 new_inactive 274 3826 3831 3842 4006 4922 4923
0520 new_inactive 273 1124 2641 2840 3849 3868 3869 4083 4334 4374 4942 4944
0521 new_inactive 396 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0522 new_inactive 370 4342 4360 6317 6348
0523 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0524 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0525 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0526 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0527 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0528 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0529 new_inactive 370 4342 4360 6317 6348
0530 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0531 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0532 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0533 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0534 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0535 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0536 new_inactive 370 4342 4360 6317 6348
0537 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0538 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0539 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0540 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0541 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0542 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0543 new_inactive 370 4342 4360 6317 6348
0544 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0545 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0546 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0547 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0548 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0549 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0550 new_inactive 370 4342 4360 6317 6348
0551 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0552 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0553 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0554 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0555 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0556 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0557 new_inactive 370 4342 4360 6317 6348
0558 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0559 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0560 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0561 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0562 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0563 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0564 new_inactive 370 4342 4360 6317 6348
0565 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0566 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0567 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0568 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0569 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0570 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0571 new_inactive 370 4342 4360 6317 6348
0572 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0573 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0574 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0575 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0576 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0577 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0578 new_inactive 370 4342 4360 6317 6348
0579 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0580 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0581 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0582 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0583 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0584 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0585 new_inactive 370 4342 4360 6317 6348
0586 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0587 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0588 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0589 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0590 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0591 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0592 new_inactive 370 4342 4360 6317 6348
0593 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0594 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0595 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0596 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0597 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0598 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0599 new_inactive 370 4342 4360 6317 6348
0600 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0601 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0602 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0603 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0604 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0605 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0606 new_inactive 370 4342 4360 6317 6348
0607 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462
0608 new_inactive 323 1281 1282 1316 1317 1333 1334 1342 1343 1344 1380 1381 1388 1397
0609 new_inactive 367 1666 1661 2214 2223 2244 2246 2498 2610 3483 6397
0610 new_inactive 368 1601 1643 1646 2188 2193 3126 4382 4688 6097 6136 6398
0611 new_inactive 291 369 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0612 new_inactive 368 1316 1333 1360 1386 1647 4490 4498 4617 6212 6241 6248 6291
0613 new_inactive 370 4342 4360 6317 6348
0614 new_inactive 371 1468 1528 1722 1728 2074 3134 4439 4461 4462

```

0009          5318
0010      total_tr      277 300 3946 3940 4366 5303 5370 5801
0011      total_u_bus   301 3047 5304 5802
0012
0013      total_u_cpu   303 3047 5306 5804
0014
0015      total_u_disk  302 3047 5306 5803
0016
0017      total_tr_time 300 3940 5303 5370 5801
0018
0019      tr_of         394 482 481 483 485 1016 1111 1209 1266 1446 1668 2083 2168 2264 2418
0020
0021      tr_phase     3123 4081 4084 4281 4354 4444 4810 4882 4921
0022
0023      tr_time      396 472 2768 2907 3300 3423 3901 4372 5077 5352
0024
0025      tr_time_std  396 1970 3640 4408 4410 4412 4418 4420 4422
0026
0027      trace_time   324
0028
0029      trace_tr     397
0030
0031      updated_jacc 390 1416 1421 3990 4622 4708 6402
0032
0033      updated_jacc 394 1423 3982 4810 4716 6396
0034
0035      updated_rsl  470 3080 4319 4322 4814 4702 4986 5303
0036
0037      working_time 282 3716
0038
0039      working_rsl  401 474 476 1410 2023 2802 3076 3090 3289 3340 3348 3349 3366 3387
0040
0041      writing_rsl   356 3569 3584 3774 3778 4608 4703 6009
0042
0043      */

```

VITA

Yao-Nan Lien was born in Taichung, Taiwan, Republic of China, on February 15, 1951. He graduated from The First N.C.O School of Chinese Army in January 1969, and served as a sergeant in Chinese Army from 1969 to 1975. He began his college education at the National Cheng Kung University in September 1975, and received a B.S. degree in Electrical Engineering in June 1979. From July 1979 to July 1980, he was a teaching assistant at the National Taiwan Institute of Technology.

He attended Purdue University for his graduate study in August 1980. Since then he had been a research assistant in the School of Electrical Engineering. He received a M.S. degree in Electrical Engineering in December 1981. He is an instructor in the Department of Computer and Information Science, the Ohio State University, since January 1986.

He is a student member of the IEEE's computer society, ACM, and AAAI.