

CERTIFICATE

It is certified that the work contained in the thesis titled “Design of Adaptive Fault-Tolerant Routing Algorithms for Hypercubes” by “Lokendra Singh Umrao” has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

It is further certified that the student has fulfilled all the requirements of Comprehensive, Candidacy and SOTA.

Dr. Ravi Shankar Singh

Department of Computer Science and Engineering

Indian Institute of Technology (BHU)

DECLARATION BY THE CANDIDATE

I, Lokendra Singh Umrao, certify that the work embodied in this thesis is my own bona fied work and carried out by me under the supervision of Dr. Ravi Shankar Singh from July 2012 to November 2015 at the Department of Computer Science and Engineering, Indian Institute of Technology, Varanasi. The matter embodied in this thesis has not been submitted for the award of any degree/diploma. I declare that I have faithfully acknowledged and given credits to the research workers wherever their works have been cited in my work in this thesis. I further declare that I have not willfully lifted up any other's work, paragraphs, text, data, results, *etc.*, reported in journals, books, magazines, reports dissertations, theses, *etc.*, or available at websites and included them in this thesis and cited as my own work.

Date : 06/11/2015

Place : Varanasi

LOKENDRA SINGH UMRAO

CERTIFICATE BY THE SUPERVISOR

It is certified that the above statement made by the student is correct to the best of my knowledge.

Dr. Ravi Shankar Singh
(Supervisor)

Dr. K. K. Shukla
(Professor and Head)
Department of Computer Science and Engineering

COPYRIGHT TRANSFER CERTIFICATE

Title of the Thesis : Design of Adaptive Fault-Tolerant Routing Algorithms for Hypercubes

Name of the Student : Lokendra Singh Umrao

Copyright Transfer

The undersigned hereby assigns to the Indian Institute of Technology (Banaras Hindu University) Varanasi all rights under copyright that may exist in and for the above thesis submitted for the award of the Ph.D.

Date : 06/11/2015

Place : Varanasi

LOKENDRA SINGH UMRAO

Note: However, the author may reproduce or authorize others to reproduce materials extracted verbatim from the thesis or derivative of the thesis for author's personal use provided that the source and the Institute's copyright notice are indicated.

Preface

The aim of this thesis is to study, implement and evaluate fault-tolerant routing algorithms for the hypercube interconnection network. We have addressed and designed fault tolerant routing algorithms in the presence of high number of node and/or link faults. We addressed this issue by designing adaptive routing protocols for hypercube interconnection networks. This technique addresses network latency and bandwidth utilization for parallel architectures. Adaptive routing algorithms exploit gains of path redundancy in n -cube.

Interconnection networks play an important role in the performance of modern high performance computing systems. It consists of a series of nodes and links. Nodes interact with each other for communication through links. The interconnection network is a requirement of any parallel computer as it helps parallel systems in showing high performance by providing reliable and quick communication over the networks. Since most routing algorithms for parallel computers are not being designed to tolerate faults, component failures impact these systems terribly. Thus, one link and/or node failure may halt the entire computing system altogether and stop the scientific applications running on them.

In this thesis, we present fault-tolerant routing algorithms based on adaptive protocols. Adaptive routing protocols can use alternative paths between communicating nodes. Multipath networks and adaptive routing protocols dynamically adapt to network conditions, thus capable of serving interconnection networks affected by a large number of node/link failures. Three contributions are presented throughout this thesis, namely: fault-tolerant distributed node-to-node routing, fault-tolerant node-to-set disjoint-path routing, and reliable broadcasting via independent spanning trees.

The aim of this thesis is to further study parallel computing, interconnection networks, routing, fault tolerance and node-disjoint paths. The scope of research given in this thesis is to design, implement and evaluate fault-tolerant routing strategies for the hypercube topology which can be utilized to frame the supercomputers.

Chapter 1 introduces parallel computing. It deals with the need for high performance computing and discusses the hypercube interconnection networks. General introduction to other related topics in parallel computing is covered to provide the span of the field.

Chapter 2 explains the basic terminologies of fault tolerance, then background on interconnection networks for HPC systems including topologies and routing. Then the concepts about hypercube interconnection networks and their routing methods and previous related works that explain routing strategy for both cases, routing without failures and routing in the presence of failures. The chapter provides complete theoretical and practical implementation of designing simple routing algorithms.

Chapter 3 describes in detail the first adaptive fault-tolerant node-to-node routing algorithm over all shortest node-disjoint paths in n -dimensional hypercube interconnection networks. It is designed in such a way that it can handle large number of node and link failures, while delivering all n messages over disjoint-paths in the presence of maximum permissible node/link failures. The chapter develops the idea from nodes and /or links failures in hyper networks towards fault tolerating interconnection networks.

Chapter 4 presents a node-to-set node-disjoint fault-tolerant routing algorithm based on subcubes of the hypercube networks. The n -dimensional hypercube can tolerate maximum $n-1$ faulty nodes. The proposed algorithm generates node disjoint-paths which maximise the probability of setting up non-faulty path in a faulty

environment.

Chapter 5 introduces data broadcasting on parallel computers through multiple independent spanning trees (ISTs). The n -IST based broadcasting from common root r on the hypercube network can provide n -degree fault tolerance. The designed fault-tolerant broadcasting algorithm using ISTs may increase message security in hypercube network.

Chapter 6 concludes the thesis and presents future directions in the research. This chapter also gives initiation to a broad range of open lines for fault-tolerant routing and further work.

In this thesis, we have designed, implemented and evaluated the different algorithms for fault-tolerant routing. The methodologies of all the algorithms are based on existing theories, knowledge and observations. The algorithms proposed in this thesis were developed on a theoretical basis and were implemented practically. With the help of relevant books and related research papers, we have focussed on the design, implementation, and evaluation of similar algorithms with efficient complexity. We have analysed the effectiveness of all the proposed fault-tolerant routing algorithms through simulation. For this, we have developed simulation models for experimental evaluation of our propositions.

Most of the concepts in the thesis are illustrated by several examples. This thesis can be useful to students and engineers who are interested in routing algorithms of high speed interconnection networks.

November, 2015
IIT (BHU), Varanasi, India

Lokendra Singh Umrao

Abstract

Interconnection networks play an important role for the performance of modern high performance computing systems. It consists of a series of nodes and links. Each node interacts with each other for communications through links. The interconnection network is a requirement of any parallel computer because parallel system shows high performance by providing reliable and quick communication over the networks. In this context, components failures have an extremely high impact because most of the routing algorithms have not been designed to tolerate faults. Since, one link and/or node failure may halt the entire computing system and stop the scientific applications running on them.

In this thesis, we present fault-tolerant routing algorithms based on adaptive protocols. Adaptive routing protocols can use alternative paths between communicating nodes. Multipath networks and adaptive routing protocols dynamically adapt to network conditions, thus capable of serving interconnection networks affected by a large number of node/link failures. Three contributions are presented throughout this thesis, namely: fault-tolerant distributed node-to-node routing, fault-tolerant node-to-set disjoint-path routing, and reliable broadcasting via independent spanning trees.

The aim of this thesis is to study, implement and evaluate fault-tolerant routing algorithms for the hypercube interconnection network. We have addressed and designed fault tolerant routing algorithms in the presence of high number of node and/or link faults. We addressed this issue by designing adaptive routing protocols for hypercube interconnection networks. This technique addresses network latency and bandwidth utilization for parallel architectures. Adaptive routing algorithms exploit gains of path redundancy in n -cube.

The first contribution of this thesis is the adaptive fault-tolerant routing algorithm for hypercube topology. This algorithm has been designed in such a way that they can use alternative path available in hypercube topology, making more efficient use of network bandwidth and allowing hypercube networks to perform in the presence of large number of faults. The proposed algorithm is a simple uniform distributed algorithm that can tolerate a large number of process failures, while delivering all n messages over optimal-length disjoint paths. However, no distributed algorithm uses acknowledgement messages (*acks*) for fault tolerance. So, for dealing the faults, acknowledgement messages (*acks*) are included in the proposed algorithm for routing messages over node-disjoint paths in the hypercube network. Simulation results confirm that the proposed node-to-node routing algorithm provides an average of 10% improvement in the performance of hypercube network in comparison with the previously proposed routing algorithms—depth first search algorithm and unsafety vectors algorithm.

The second contribution is the node-to-set node-disjoint routing algorithm for the hypercube networks with faulty nodes. This algorithm has been designed to the problems of the disjoint shortest paths routing. The proposed algorithm can tolerate maximum $n - 1$ faulty nodes, where n is the dimension of the hypercube. The proposed NoSeRo algorithm used the subcube property of the n -dimensional hypercube. It adapts divide-and-conquer approach to take full advantage of the regularity of the hypercube. Hence, proposed algorithm generates fault-free node-disjoint paths in a faulty environment. The proposed fault tolerant routing algorithm for faulty hypercube networks which finds n disjoint paths from source process s to n destination processes in n -dimensional hypercube in $O(n^2)$ time with optimal path lengths at most $n + f + 1$, where n is the number of destination node and f is the number of faulty nodes. Then simulation results showed that the proposed algorithm reduce the average path length by about 20% in comparison of Bossard's algorithm in 8-dimensional hypercube (H_8).

The third contribution is the reliable data broadcasting scheme by generating Independent Spanning Trees (ISTs) on hypercubes. The proposed

scheme can be useful for secure message transmission. Using n -IST-based broadcasting from same root r on hypercube network ($N = 2^n$) provides n -degree fault tolerance. The proposed algorithm can be easily implemented in parallel or distributed systems. Using ISTs one can enhance the fault-tolerance, bandwidth, and security. In this chapter, we study the existence and construction of n ISTs rooted at an arbitrary vertex in H_n ($n \geq 1$). A parallel algorithm with the time complexity $O(n)$ is proposed to construct n ISTs on H_n , where $n \geq 1$.

In this thesis, we have designed, implemented and evaluated the different algorithms for fault-tolerant routing. The methodologies of all the algorithms are based on existing theories, knowledge and observations. The algorithms proposed in this thesis were developed on a theoretical basis and were implemented practically. With the help of relevant books and related research papers, we have focussed on the design, implementation, and evaluation of similar algorithms with efficient complexity. We have analysed the effectiveness of all the proposed fault-tolerant routing algorithms through simulation. For this, we have developed simulation models for experimental evaluation of our propositions.

All of the proposals made in this thesis are suitable (without hardware modification) to be implemented on currently personal computing systems and all the proposed algorithms are able to tolerate dynamically a reasonable number of faults.

Keywords: Hypercube interconnection networks, Fault tolerance, Fault-tolerant routing, Node-disjoint paths, Independent spanning trees, Multicasting, Broadcasting.

To my family. They give me the solid basis that always me to reach here.

and

especially to my wife Madhulika and my daughter Pihu.

Acknowledgements

First of all, I would like to express my profound gratitude to my guide, Dr. Ravi Shankar Singh. His guidance and patience made this dissertation possible. His instruction and insight from the formative stages of the research through the final stages provided valuable direction to my work.

I am also thankful to the members of my committee, Dr. K. K. Shukla, Dr. A. K. Tripathi, and Dr. Ashokji Gupta for their encouragement and advice on my research, as well as Dr. S. K. Singh who served as the DPGC Representative on my committee.

Special thanks are due to Dr. R. B. Mishra, Dr. Rajeev Srivastava, Dr. Vinayak Srivastava, Dr. Bhaskar Biswas and Dr. Ravindranath Chowdary C, who have given me their constant support and advice. I am also grateful to my friends, Subhash Chandra Patel, Dharmendra Prasad Mahato, Sumit Jaiswal and Ali Imam Abidi. As a friends, they continuously encouraged and supported me in many ways and made my graduate life more enjoyable and memorable.

Most of all, thanks to my family who were always there when I needed them. My parents never lost trust in me and always prayed to God for me. They were always supportive throughout my research and gave me unconditional love. I dedicate this dissertation to my family.

Contents

List of Figures	vii
List of Tables	xi
1 Introduction	1
1.1 Parallel Computing	1
1.1.1 Interconnection Network	5
1.1.2 Fault Tolerance in Network	5
1.2 Problem Statement	6
1.3 Motivation	6
1.4 Objectives	7
1.5 Contributions	8
1.6 Research Methodologies	9
1.7 Scope and Organization of the Thesis	10
2 Thesis Background	11
2.1 Dependability in Distributed Systems	11
2.1.1 Distributed Systems Terminology	11
2.1.2 Dependable High Performance Computing Systems	13
2.2 Fault Tolerance	13
2.2.1 Fault Tolerance Terminology	13
2.2.2 Fault-Tolerant Routing	16
2.3 Interconnection Networks	18
2.3.1 Network Topologies	19
2.3.2 Routing	21
2.3.2.1 Deterministic Routing	22

CONTENTS

2.3.2.2	Adaptive Routing	23
2.4	Hypercube Basics	25
2.4.1	Subcube Reliability Computation in Hypercube Networks	27
2.5	Fault Tolerance for Hypercube Interconnection Networks	28
2.5.1	Hardware-based Fault Tolerance in Hypercube Networks	28
2.5.2	Software-based Fault Tolerance in Hypercube Networks	28
2.5.2.1	Hypercube Free Dimension Technique	28
2.5.2.2	Spanning Tree Approach	29
2.5.3	Fault-Tolerant Routing Techniques in Hypercube Networks	30
2.5.3.1	Depth-First Search based Routing	31
2.5.3.2	Heuristic-based Routing	33
2.5.3.3	Node Safety based Routing	34
2.5.3.4	Reachability based Routing	35
2.5.4	Multicasting/Broadcasting in Faulty Hypercube	36
3	Fault-Tolerant Distributed Node-to-Node Routing	39
3.1	Overview	40
3.2	Fault-Tolerant Distributed Routing	42
3.2.1	Detection & Notification	42
3.2.2	Selection of Alternative Paths	43
3.2.3	Configuration of Alternative Paths	44
3.2.4	Transient and Permanent Faults	45
3.3	Preliminaries	48
3.3.1	Basic Hypercube Routing	49
3.4	Algorithms	51
3.4.1	Routing Without Failures	53
3.4.2	Routing With Failures	54
3.5	Proof of Correctness	56
3.6	Evaluation Model	57
3.7	Implementation & Simulation Results	63
3.8	Conclusion	70

4	Fault-Tolerant Node-to-Set Disjoint-Path Routing	73
4.1	Overview	74
4.2	Subcube Reliability Computation	76
4.3	Algorithm	76
4.3.1	Fault-Tolerant Node-to-Node Routing in Hypercube	78
4.3.2	The Proposed Node-to-Set Routing Algorithm	79
4.4	Complexities Analysis	81
4.5	Evaluation Model	83
4.6	Performance Evaluation	83
4.7	Conclusion	88
5	Reliable Broadcasting via Independent Spanning Trees	89
5.1	Overview	90
5.1.1	Parent Exchange-based n-IST Optimal Construction	93
5.1.2	HDLS-based n-IST Parallel Construction	93
5.2	Constructing Independent Spanning Trees on Hypercube	95
5.2.1	Virtual Roots for Reliable Broadcasting	96
5.2.2	Parallel Preprocessing for Generating n ISTs	96
5.2.3	Algorithm	97
5.3	Correctness	100
5.4	Experimental Results	105
5.4.1	Response Time of the Parallel Construction of n ISTs	105
5.4.2	Reliable Broadcasting	106
5.5	Conclusion	107
6	Summary and Conclusions	109
6.1	Conclusions	110
6.2	Scope for Further Work	111
6.3	List of Publications	111
6.4	List of Accepted Papers	112
6.5	List of Communicated Papers	112
	Bibliography	113

CONTENTS

List of Figures

1.1	Steps of solving a problem using parallel computing.	3
2.1	The dependability tree.	12
2.2	The relationship among fault, error and failure.	14
2.3	Network failure terminology.	15
2.4	Classification of direct network topologies.	20
2.5	Taxonomy of routing algorithms.	21
2.6	XY routing in a mesh with $N = 16$. Messages are routed in a dimension-ordered fashion.	23
2.7	e-cube routing in a hypercube with $N = 16$. Messages are routed in a dimension-ordered fashion.	24
2.8	A taxonomy for adaptive routing protocols.	24
2.9	A 4-dimensional hypercube and its subcubes.	26
2.10	The structure of the dimension in the space.	29
2.11	Hypercube division using free dimension.	29
2.12	A Spanning tree (ST) of the hypercube: First phase.	30
2.13	The modified ST under faulty node: Second reconfiguration phase.	30
2.14	Routing in 4-cube using DFS.	32
2.15	Shortest path routing in 4-cube.	33
2.16	Example of node status in 4-cube.	34
2.17	Example of strongly and ordinary unsafe nodes in 4-cube.	35
2.18	Example of node reachability in 4-cube.	36
2.19	Multicast tree in H_5 with reference to Table 2.2.	37
3.1	Message passing in normal condition.	42

LIST OF FIGURES

3.2	Notification of node/link failures.	42
3.3	Selection of alternative path.	43
3.4	A multipath composed by 4 multistep paths in hypercube.	45
3.5	Reception of <i>acks</i> messages.	47
3.6	Message passing flow diagram.	48
3.7	CPN model of the system.	58
3.8	CPN model of the system after few steps.	61
3.9	CPN model of the system after all steps.	61
3.10	CPN model for the Sender Side.	62
3.11	CPN model of the Receiver acks.	62
3.12	CPN model of the Transmit Message.	63
3.13	CPN model of the Transmit acks.	63
3.14	CPN model of the receiver.	64
3.15	Simulation result for different message lengths with 30% faulty nodes in H_8	65
3.16	Simulation result for different message lengths with 40% faulty nodes in H_8	66
3.17	Simulation result for different message lengths with 50% faulty nodes in H_8	66
3.18	Simulation result for different message lengths with 60% faulty nodes in H_8	67
3.19	Simulation results for different generation rates with 30% faulty nodes in H_8	67
3.20	Simulation results for different generation rates with 40% faulty nodes in H_8	68
3.21	Simulation results for different generation rates with 50% faulty nodes in H_8	68
3.22	Simulation results for different generation rates with 60% faulty nodes in H_8	69
4.1	Embedding is not possible in a faulty 4-cube since no fault-free 3-cube exist.	77

LIST OF FIGURES

4.2	Embedding is not possible in a faulty 4-cube since no fault-free 3-cube exist.	77
4.3	Embedding is possible in a faulty 4-cube since fault-free 3-cube exists. .	77
4.4	Embedding is possible in a faulty 4-cube since fault-free 3-cube exists. .	78
4.5	A 5-dimensional hypercube.	81
4.6	CPN model for computing.	84
4.7	CPN model for date generation.	84
4.8	CPN model for task scheduler.	85
4.9	CPN model for computing.	85
4.10	CPN model for computing.	85
4.11	CPN model for computing.	86
4.12	Comparison in terms of path length with failures for node-to-set disjoint path routing in H_8	87
4.13	Comparison in terms of path length without failures for node-to-set disjoint path routing in H_8	87
5.1	The spanning binomial trees with $r = 0$ on H_4	92
5.2	Four pairwise independent spanning trees of H_4 rooted at the vertex 0. .	95
5.3	Response time of the parallel construction of n ISTs	106
5.4	Ratio b/w faulty links vs total links in reliable broadcasting.	108

LIST OF FIGURES

List of Tables

2.1	Results of heuristic routing algorithm.	34
2.2	The resulting relative addresses with respect to source node 12 (01100) .	37
4.1	Comparison in terms of path length and time complexity	88
5.1	Parent of a node $x \in H_4$ in $T_0(r = 0), i = 0$	99
5.2	Parent of a node $x \in H_4$ in $T_1(r = 0), i = 1$	100
5.3	Parent of a node $x \in H_4$ in $T_2(r = 0), i = 2$	101
5.4	Parent of a node $x \in H_4$ in $T_3(r = 0), i = 3$	102
5.5	Parent of a node $x \in H_5$ in $T_4(r = 8), i = 4$	103
5.6	Comparison in terms of no. of computation steps on H_8	106
5.7	Reliable broadcasting in faulty links	107
5.8	Comparison in terms of approach and advantages	108

LIST OF TABLES

List of Algorithms

1	Algorithm for node-to-node routing in a hypercube	54
2	Algorithm for node-to-node routing in a faulty hypercube	55
3	NoSeRo ($H_n, D = \{d_1, d_2, \dots, d_n\}, F = \{f_1, f_2, \dots, f_{n-1}\}$)	80
4	GEN-PARENTS	94
5	Efficient parallel ISTs	97

LIST OF ALGORITHMS
