

Do-All Computing in Distributed Systems

*Cooperation in the
Presence of Adversity*

Do-All Computing in Distributed Systems

Cooperation in the Presence of Adversity

by

Chryssis Georgiou
*University of Cyprus
Cyprus*

and

Alexander A. Shvartsman
*Massachusetts Institute of Technology (MIT)
USA*

 Springer

Chryssis Georgiou
University of Cyprus
Dept. Computer Science
P.O.Box 20537
1678 Nicosia, CYPRUS
chryssis@ucy.ac.cy

Alexander A. Shvartsman
University of Connecticut
Computer Science and Engineering
371 Fairfield Way
Storrs, CT 06268, USA
aas@cse.uconn.edu

Library of Congress Control Number: 2007937388

Do-All Computing in Distributed Systems: Cooperation in the Presence of Adversity
by Chryssis Georgiou and Alexander A. Shvartsman

ISBN-13: 978-0-387-30918-7

e-ISBN-13: 978-0-387-69045-2

Printed on acid-free paper.

© 2008 Springer Science+Business Media, LLC.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

9 8 7 6 5 4 3 2 1

springer.com

To my wife Agni, and son Yiorgo
CG

To Sana, my wife and best friend
AAS

Contents

List of Figures	XI
List of Symbols	XIII
Foreword by Michel Raynal	XV
Authors' Preface	XIX
1 Introduction	1
1.1 Do-All Computing	2
1.2 Do-All and Adversity	4
1.3 Solving Do-All: Fault-Tolerance with Efficiency	6
1.4 Chapter Notes	8
2 Distributed Cooperation Problems:	
Models and Definitions	11
2.1 Model of Computation	11
2.1.1 Distributed Setting	11
2.1.2 Communication	11
2.2 Models of Adversity	12
2.2.1 Processor Failure Types	12
2.2.2 Network Partitions	13
2.2.3 Adversaries and their Behavior	13
2.3 Tasks and Do-All Computing	14
2.3.1 The Do-All Problem	15
2.3.2 The Omni-Do Problem	16
2.4 Measures of Efficiency	17
2.5 Chapter Notes	19

3	Synchronous Do-All with Crashes: Using Perfect Knowledge and Reliable Multicast	21
3.1	Adversarial Model	22
3.2	Lower and Upper Bounds for Abstract Models	22
3.2.1	Modeling Knowledge	22
3.2.2	Lower Bounds	23
3.2.3	Upper Bounds	28
3.3	Solving Do-All Using Reliable Multicast	33
3.3.1	Algorithm AN	34
3.3.2	Correctness of algorithm AN	38
3.3.3	Analysis of Algorithm AN	40
3.3.4	Analysis of Message-Passing Iterative Do-All	44
3.4	Open Problems	45
3.5	Chapter Notes	45
4	Synchronous Do-All with Crashes and Point-to-Point Messaging	47
4.1	The Gossip Problem	48
4.2	Combinatorial Tools	49
4.2.1	Communication Graphs	49
4.2.2	Sets of Permutations	50
4.3	The Gossip Algorithm	51
4.3.1	Description of Algorithm GOSSIP _{ϵ}	51
4.3.2	Correctness of Algorithm GOSSIP _{ϵ}	55
4.3.3	Analysis of Algorithm GOSSIP _{ϵ}	59
4.4	The Do-All Algorithm	61
4.4.1	Description of Algorithm DOALL _{ϵ}	62
4.4.2	Correctness of Algorithm DOALL _{ϵ}	64
4.4.3	Analysis of Algorithm DOALL _{ϵ}	67
4.5	Open Problems	72
4.6	Chapter Notes	73
5	Synchronous Do-All with Crashes and Restarts	77
5.1	Adversarial Model	78
5.2	A Lower Bound on Work for Restartable Processors	79
5.3	Algorithm AR for Restartable Processors	82
5.3.1	Description of Algorithm AR	82
5.3.2	Correctness of Algorithm AR	86
5.3.3	Complexity Analysis of Algorithm AR	89
5.4	Open Problems	92
5.5	Chapter Notes	93

6	Synchronous Do-All with Byzantine Failures	95
6.1	Adversarial Model	96
6.2	Task Execution without Verification	96
6.2.1	Known Maximum Number of Failures	96
6.2.2	Unknown Maximum Number of Failures	98
6.3	Task Execution with Verification	98
6.3.1	Known Maximum Number of Failures	99
6.3.2	Unknown Maximum Number of Failures	111
6.4	Open Problems	112
6.5	Chapter Notes	112
7	Asynchrony and Delay-Sensitive Bounds	115
7.1	Adversarial Model and Complexity	116
7.2	Delay-Sensitive Lower Bounds on Work	118
7.2.1	Deterministic Delay-Sensitive Lower Bound	119
7.2.2	Delay-sensitive Lower Bound for Randomized Algorithms	121
7.3	Contention of Permutations	125
7.3.1	Contention and Oblivious Tasks Scheduling	127
7.3.2	Generalized Contention	128
7.4	Deterministic Algorithms Family DA	130
7.4.1	Construction and Correctness of Algorithm DA(q)	131
7.4.2	Complexity Analysis of Algorithm DA(q)	134
7.5	Permutation Algorithms Family PA	137
7.5.1	Algorithm Specification	137
7.5.2	Complexity Analysis	139
7.6	Open Problems	142
7.7	Chapter Notes	143
8	Analysis of Omni-Do in Asynchronous Partitionable Networks	145
8.1	Models of Adversity	146
8.2	A Group Communication Service and Notation	148
8.3	View-Graphs	150
8.4	Algorithm AX	154
8.4.1	Description of the Algorithm	154
8.4.2	Correctness of the Algorithm	155
8.5	Analysis of Algorithm AX	158
8.5.1	Work Complexity	158
8.5.2	Message Complexity	162
8.5.3	Analysis Under Adversary \mathcal{A}_F	165
8.6	Open Problems	165
8.7	Chapter Notes	166

9	Competitive Analysis of Omni-Do in Partitionable Networks	169
9.1	Model of Adversity, Competitiveness and Definitions	170
9.1.1	Adversary \mathcal{A}_{GR}	171
9.1.2	Measuring Competitiveness	173
9.1.3	Formalizing Computation Width	174
9.2	Algorithm RS and its Analysis	175
9.2.1	Description of Algorithm RS	175
9.2.2	Analysis of Algorithm RS	175
9.2.3	Deterministic Algorithms	178
9.3	Lower Bounds	179
9.4	Open Problems	181
9.5	Chapter Notes	181
10	Cooperation in the Absence of Communication	183
10.1	Adversity, Schedules, Waste, and Designs	184
10.2	Redundancy without Communication: a Lower Bound	187
10.3	Random Schedules	188
10.4	Derandomization via Finite Geometries	190
10.5	Open Problems	192
10.6	Chapter Notes	192
11	Related Cooperation Problems and Models	195
11.1	<i>Do-All</i> in Shared-Memory	195
11.2	<i>Do-All</i> with Broadcast Channels	200
11.3	Consensus and its Connection to <i>Do-All</i>	202
	References	205
	Index	213

List of Figures

3.1	Oracle-based algorithm.	28
3.2	A local view for phase $\ell + 2$	35
3.3	Phase ℓ of algorithm AN	37
3.4	A phase of algorithm AN.	37
4.1	Algorithm GOSSIP $_{\varepsilon}$, stated for processor v ; $\pi_v(i)$ denotes the i^{th} element of permutation π_v	52
4.2	Algorithm DOALL $_{\varepsilon}$, stated for processor v ; $\pi_v(i)$ denotes the i^{th} element of permutation π_v	63
4.3	Classification of a phase i of epoch ℓ ; execution ξ is implied. . . .	67
5.1	A local view for phase $\ell + 4$	83
5.2	Phase ℓ of algorithm AR (text in <i>italics</i> highlights differences between algorithm AR and algorithm AN).	84
5.3	A phase of algorithm AR.	85
6.1	Algorithm <i>Cover</i> . The code is for processor q	97
6.2	Algorithm for the case $f \geq p/2$. The code is for processor q . The call to the procedure is made with $P = [p]$, $T = [n]$, and $\psi = f$	101
6.3	Algorithm for the case $f < p/2$. The code is for processor q . The call parameters are $P = [p]$, $T = [n]$, and $\psi = f$	105
6.4	Subroutine <i>Do_Work_and_Verify</i> . Code for processor q	107
6.5	Subroutine <i>Checkpoint</i> . Code for processor q	108
7.1	Illustration of the adversarial strategy leading to the delay-sensitive lower bound on total-work for randomized algorithms.	125
7.2	Algorithm OBLIDO.	127
7.3	The deterministic algorithm DA ($p \geq n$).	132

7.4 Permutation algorithm and its specializations for PARAN1,
 PARAN2, and PADET ($p \geq n$). 138

8.1 Example of a view-graph 151

8.2 Input/Output Automata specification of algorithm AX. 156

9.1 An example of a $(12, n)$ -DAG. 172

List of Symbols

p	number of processors	2
n	number of tasks	2
PID or pid	unique processor identifier	3
\log	logarithm to the base 2	5
\mathcal{P}	set of processor ids numbered from 1 to p	11
\mathcal{A}	adversary or adversarial model	14
A	an algorithm	14
$\mathcal{E}(A, \mathcal{A})$	set of all executions of algorithm A for adversary \mathcal{A}	14
ξ	an execution in $\mathcal{E}(A, \mathcal{A})$	14
$\xi _{\mathcal{A}}$	adversarial pattern of ξ caused by \mathcal{A}	14
$\ \xi _{\mathcal{A}}\ $	the weight of adversarial pattern $\xi _{\mathcal{A}}$	14
TID or tid	unique task identifier	15
\mathcal{T}	set of task ids numbered from 1 to n	15
$Do\text{-}All_{\mathcal{A}}(n, p, f)$	<i>Do-All</i> problem for n tasks, p processors and adversary \mathcal{A} constrained to adversarial patterns of weight less or equal to f	15
$r\text{-}Do\text{-}All_{\mathcal{A}}(n, p, f)$	<i>iterative Do-All</i> problem for r sets of n tasks, p processors, and adversary \mathcal{A} constrained to adversarial patterns of weight less or equal to f	15
$Omni\text{-}Do_{\mathcal{A}}(n, p, f)$	<i>Omni-Do</i> problem for n tasks, p processors and adversary \mathcal{A} constrained to adversarial patterns of weight less or equal to f	16
S or $S_{\mathcal{A}}(n, p, f)$	total-work or available processor steps complexity required for an algorithm to solve a problem of size n using p processors under adversary \mathcal{A} restricted to adversarial patterns of weight no more than f	17
$ES_{\mathcal{A}}(n, p, f)$	Expected total-work complexity	17
W or $W_{\mathcal{A}}(n, p, f)$	task-oriented work complexity	17
$EW_{\mathcal{A}}(n, p, f)$	Expected task-oriented work complexity	18
M or $M_{\mathcal{A}}(n, p, f)$	message complexity	18

$EM_{\mathcal{A}}(n, p, f)$	Expected message complexity	19
\mathcal{A}_C	adversary causing processor crashes	22
$Do\text{-}All_{\mathcal{A}}^{\mathcal{O}}(n, p, f)$	$Do\text{-}All_{\mathcal{A}}(n, p, f)$ problem where processors are assisted by oracle \mathcal{O}	23
$r\text{-}Do\text{-}All_{\mathcal{A}}^{\mathcal{O}}(n, p, f)$	$r\text{-}Do\text{-}All_{\mathcal{A}}(n, p, f)$ problem where processors are assisted by oracle \mathcal{O}	23
$Gossip_{\mathcal{A}}(p, f)$	$Gossip$ problem for p processors and adversary \mathcal{A} constrained to adversarial patterns of weight less or equal to f	48
\mathcal{S}_t	symmetric group	50
\mathcal{A}_{CR}	adversary causing processor crashes and restarts . .	78
$\mathcal{A}_{CR}^{(\kappa)}$	maximal subset of \mathcal{A}_{CR} that contains only κ -restricted adversarial patterns	79
\mathcal{A}_B	adversary causing Byzantine processor failures . . .	96
\mathcal{A}_D	adversary causing arbitrary delays between local processors steps and arbitrary message delays . . .	116
$\mathcal{A}_D^{(d)}$	adversary causing arbitrary delays between local processors steps and message delays up to d time units	116
\mathcal{A}_F	adversary causing group fragmentations	146
\mathcal{A}_{FM}	adversary causing group fragmentations and merges	147
\mathcal{A}_{GR}	adversary causing arbitrary regroupings	171
\mathcal{A}_R	adversary causing rendezvous	185
$\mathcal{A}_R^{(r)}$	adversary causing at most r -way rendezvous	185

Foreword

Distributed computing was born in the late 1970s when researchers and practitioners started taking into account the intrinsic characteristics of physically distributed systems. The field then emerged as a specialized research area distinct from networking, operating systems, and parallel computing. *Distributed computing* arises when one has to solve a problem in terms of distributed entities, usually called processors, nodes, agents, sensors, peers, actors, processes, etc., such that each entity has only a partial knowledge of the many parameters involved in the problem that has to be solved. While parallel computing and real-time computing can be characterized respectively by the terms *efficiency* and *on time computing*, distributed computing can be characterized by the term *uncertainty*. This uncertainty is created by asynchrony, failures, unstable behaviors, non-monotonicity, system dynamism, mobility, connectivity instability, etc. Mastering one form or another of uncertainty is pervasive in all distributed computing problems.

The unprecedented growth of the Internet as a massive distributed network in the last decade created a platform for new distributed applications that in turn poses new challenges for distributed computing research. One such class of distributed applications is comprised of computing-intensive problems that in the past were relegated to the realm of massively parallel systems. The Internet, with its millions of interconnected computers, presents itself as a natural platform where the availability of massive distributed computing resources is seen as a compelling alternative to expensive specialized parallel supercomputers. Large networks, used as distributed supercomputers, scale much better than tightly-coupled parallel machines while providing much higher potential for parallel processing. However, harnessing the computing power contained within large networks is challenging because, unlike the applications developed for the controlled computing environments of purposefully-designed parallel systems, applications destined for distributed systems must exist in the environment fraught with uncertainty and adversity.

The field of distributed computing research, as many other areas of informatics, has traditionally encompassed both *science* and *engineering* dimensions. Roughly speaking, these can be seen as complementary facets: science is to understand and engineering is to build. With respect to distributed computing, we are often concerned with a *science of abstraction*, namely, creating the right model for a problem and devising the appropriate mechanizable techniques to solve it. This is particularly true in fault-tolerant, dynamic, large-scale distributed computing where finding models that are realistic while remaining abstract enough to be tractable, was, is, and still remains a real challenge.

The monograph by Chryssis Georgiou and Alex Shvartsman presents a very comprehensive study of massive cooperative computing in distributed settings in the presence of adversity. They focus on a problem that meaningfully abstracts a network supercomputing paradigm, specifically where distributed computing agents cooperate on performing a large number of independent tasks. Such a computation paradigm forms a cornerstone for solutions to several computation-intensive problems ranging from distributed search to distributed simulation and multi-agent collaboration. For the purposes of this study, the authors define *Do-All* as the problem of multiple processors in a network cooperatively performing a collection of independent tasks in the presence of adversity, such as processor failures, asynchrony, and breakdowns in communication. Achieving efficiency in such cooperation is difficult due to the dynamic characteristics of the distributed environments in which computing agents operate, including network failures, and processor failures that can range from the benign crash failures to the failures where faulty components may behave arbitrarily and even maliciously. The *Do-All* problem and its iterative version is used to identify the trade-offs between efficiency and fault-tolerance in distributed cooperative computing, and as a target for algorithm development. The ultimate goal is to develop algorithms that combine efficiency with fault-tolerance to the maximum extent possible, and that can serve as building blocks for network supercomputing applications and, more generally, for applications requiring distributed cooperation in the face of adversity.

During the last two decades, significant research was dedicated to studying the *Do-All* problem in various models of computation, including message-passing, partitionable networks, and shared-memory models under specific assumptions about synchrony/asynchrony and failures. This monograph presents in a coherent and rigorous manner the lower bound results and the most significant algorithmic solutions developed for *Do-All* in the message-passing model, including partitionable networks. The topics chosen for presentation include several relevant models of adversity commonly encountered in distributed computing and a variety of algorithmics illustrating important and effective techniques for solving the problem of distributed cooperation. The monograph also includes detailed complexity analysis of algorithms, assessing their efficiency in terms of work, communication, and time.

As the aim of a theory is to codify knowledge in order for it to be transmitted (to researchers, students, engineers, practitioners, etc), the research results presented in this monograph are among the fundamental bases in distributed computing theory. When effective distributed cooperation is possible, we learn why and how it works, and where there exist inherent limitations in distributed cooperation, we learn what they are and why they exist.

Rennes, France
September 2007

Michel Raynal

Authors' Preface

With the advent of ubiquitous high bandwidth Internet connections, network supercomputing is increasingly becoming a popular means for harnessing the computing power of an enormous number of processes around the world. Internet supercomputing comes at a cost substantially lower than acquiring a supercomputer or building a cluster of powerful machines. Several Internet supercomputers are in existence today, for instance, Internet PrimeNet Server, a project comprised of about 30,000 servers, PCs, and laptop computers, supported by Entropia.com, Inc., is a distributed, massively parallel mathematics research Internet supercomputer. PrimeNet Server has sustained throughput of over 1 teraflop. Another popular Internet supercomputer, the SETI@home project, also reported its speed to be in teraflops.

In such distributed supercomputing settings it is often the case that a very large number of independent tasks must be performed by an equally large number of computers. Given the massive numbers of participating computers, it is invariably the case that non-trivial subsets of these machines may be faulty, disconnected, experiencing delays, or simply off-line at any given point in time. At such scales of distributed computing, failures are no longer an exception, but the norm. For example, a visitor to the network control center at Akamai Technologies, a global Internet content and application delivery company, will immediately notice that the floor-to-ceiling monitor-paneled walls of the main control room display a surprisingly large number of server icons in red, indicating server failures. Yet the services delivered by the company's 25,000 servers worldwide continue unaffected, and there is little alarm among the engineers monitoring the displays. Dealing with failures is routine business, provided the massively distributed system has built-in redundancy and is able to combine efficiency with fault-tolerance.

In another example, Internet supercomputing, such as SETI@home, involves large sets of independent tasks performed by distributed worker computers. One of the major concerns involved in such computing environments is the reliability of the results returned by the workers. While most participating computers may be reliable, a large number of the workers have been known

to return incorrect results for various reasons. Workers may return incorrect results due to unintended failures caused, for example, by over-clocked processors, or they may claim to have performed assigned work so as to obtain incentives, such as getting higher rank on the SETI@home list of contributed units of work. This problem already exists in the setting where the task allocation is centralized, and assumed to be reliable. The problem becomes substantially more difficult when the task allocation also has to be implemented in a highly-distributed fashion to provide the much needed parallelism for computation speed-up and redundancy for fault tolerance. In such settings it is extremely important to develop distributed algorithms that can be used to ensure dependable and efficient execution of the very large numbers of tasks.

In this monograph we abstract the problem of distributed cooperation in terms of the *Do-All* problem, defined as the problem of p processors in the network, cooperatively performing n independent tasks, in the presence of adversity. In solving this problem, we pursue the goal of combining the reliability potential that comes with replicated processors in distributed computation, with the speed-up potential of performing the large number of tasks in parallel. The difficulty associated with combining fault-tolerance with efficiency is that the two have conflicting means: fault-tolerance is achieved by *introducing redundancy*, while efficiency is achieved by *removing redundancy*. We present several significant advances in algorithms designed to solve the *Do-All* problem in distributed message-passing settings under various models of adversity, such as processor crashes, asynchrony, message delays, network partitions, and malicious processor behaviors. The efficiency of algorithms for *Do-All* is most commonly assessed in terms of *work* and *communication* complexity, depending on the specific model of computation. Work is defined either as the total number of computational steps taken by all available processors during the computation or as the total number of *task-oriented* computational steps taken by the processors. A computational step taken by a processor is said to be task-oriented, if during that step the processor performs a *Do-All* task. We refer to the first variation of work as *total-work* and the second variation of work as *task-oriented work*. We develop corresponding complexity analyses that show to what extent efficiency can be combined with fault-tolerance. We also present lower bounds that capture theoretical limitations on the possibility of combining fault-tolerance and efficiency. In this work we ultimately aim to provide robust, i.e., efficient and fault-tolerant, algorithms that will help bridge the gap between abstract models of dependable network computing and realistic distributed systems.

Monograph Roadmap

In Chapter 1 we provide motivation, introduce the distributed cooperation problem *Do-All* and discuss several variants of the problem in different models of computation.

In Chapter 2 we formal the basic message-passing model of computation used in this monograph, and present several models of adversarial settings studied in subsequent chapters. We define the nature of the tasks – the input to the distributed cooperation problem. We define the *Do-All* problem, and its counterpart for partitionable networks, the *Omni-Do* problem. We conclude the chapter with the definitions of main complexity measures used in the sequel: total-work, task-oriented work, and message complexity.

In Chapter 3 we study the *Do-All* problem for distributed settings with processor crashes. We provide upper and lower bounds on work for solving *Do-All* under the assumption of perfect knowledge, e.g., when an algorithm is aided by an omniscient oracle. We put these result to use by developing an efficient and fault-tolerant algorithm for *Do-All* where processors communicate by means of reliable broadcasts.

In Chapter 4 we develop a solution for the *Do-All* problem for the setting with processor crashes, where processors communicate using point-to-point messaging. This algorithm uses a gossip algorithm as a building block, also presented in the chapter.

In Chapter 5 we give lower bounds on work for *Do-All* in the model where processors are subject to crashes and restarts, and we develop and analyze an algorithm for this model of adversity.

In Chapter 6 we study the complexity of *Do-All* in the adversarial model where processors are subject to Byzantine failures, that is, where faulty processors may behave arbitrarily and even maliciously. We provide several algorithms and lower bound results under this model of adversity.

In Chapter 7 we study the upper and lower bounds of solving *Do-All* in the setting where an adversary introduces processor asynchrony and message delays. We present several algorithm for this model and provide their delay-sensitive analysis.

In Chapter 8 we switch our attention to partitionable networks and the *Omni-Do* problem. We give an efficient algorithm that solves *Omni-Do* in the presence of network fragmentation and merges.

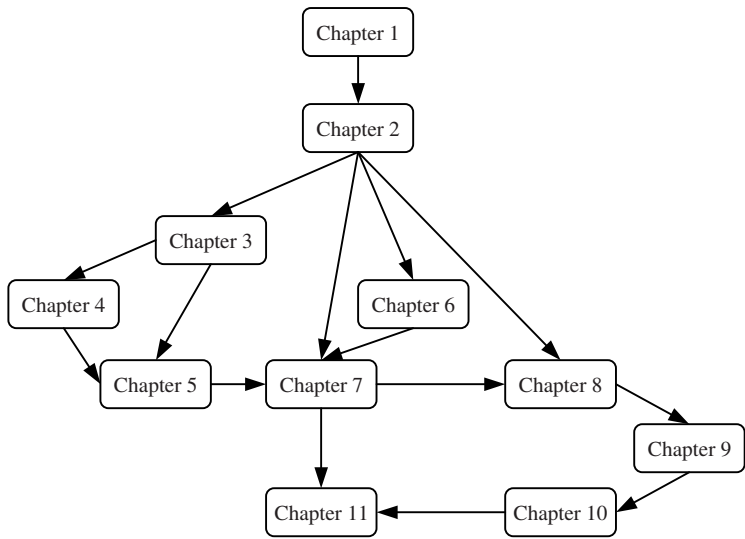
In Chapter 9 we study the *Omni-Do* problem in the model where the network can undergo arbitrary reconfigurations. We assess upper and lower bounds for the problem using competitive analysis.

In Chapter 10 we study *Do-All* in the setting where the adversary initially starts processors in isolated singleton groups, and then allows the processor to rendezvous. We analyze redundant work performed by the isolated processors prior to rendezvous, and we present several scheduling strategies designed to minimize redundant task executions.

Finally, in Chapter 11 we survey related problems and models, including the problem of distributed cooperation in shared-memory models, algorithms for the model where processors communicate through broadcast channels, and we show a connection between *Do-All* and the distributed consensus problem.

The chapters of this monograph can of course be read in the sequential order from Chapter 1 to Chapter 11. In the diagram that follows we show

alternative suggested paths through the monograph. Chapters 1 and 2 should be read in sequence before other chapters. It is also recommended that Chapters 8, 9, and 10 are read in sequence. The only remaining dependency is that Chapter 3 is read before Chapter 5.



In presenting our message-passing algorithms, we aim to illustrate the most interesting algorithmic techniques and paradigms, using a clear high-level level pseudocode that is best suited to represent the nature of each algorithm.

Each chapter concludes with an overview of open problems relevant to the topics presented in the chapter, and a section containing chapter notes, including detailed bibliographic notes, and selected comparisons with and overviews of related work.

Bibliographic Notes

At the end of each chapter we provide Chapter Notes that contain bibliographic notes and overview related topics and results. The complete bibliography follows the last chapter. Here we give additional pointers to conference proceedings, archival journals, and books covering the various areas related to distributed computing and fault-tolerant algorithms. Most results in this monograph appeared as articles in journals or conference proceedings (see bibliography), additionally the main results in Chapters 3, 4, 6, 8, and 9 appear in the PhD dissertation of the first author [43].

Work on fault-tolerant distributed computation related to the content of this monograph appear in the proceedings of conferences, in journals, and in books. A reader interested in learning more about this ongoing research as well as research beyond the scope of this volume will be well served by consulting recent publication on such topics from the venues we list below.

The following conferences are examples of the most relevant fora for results related to topics in this monograph: ACM symposium on Principles of Distributed Computing (PODC), ACM symposium on Parallel Algorithms and Architectures (SPAA), ACM symposium on Theory of Computing (STOC), ACM-SIAM symposium on Discrete Algorithms (SODA), IEEE symposium on Foundations of Computer Science (FOCS), IEEE sponsored conference on Distributed Computing Systems (ICDCS), EATCS sponsored symposium on Distributed Computing (DISC), the conference on the Principles of Distributed Systems (OPODIS) and the colloquium on Structural Information and Communication Complexity (SIROCCO). The most relevant journals include: Springer Distributed Computing, SIAM Journal on Computing, Theoretical Computer Science, Information and Computation, Information Processing Letters, Parallel Processing Letters, Journal of the ACM, Journal of Algorithms, Journal of Discrete Algorithms, and Journal of Parallel and Distributed Computing.

The 1997 book by Kanellakis and Shvartsman [67] presents research results for fault-tolerant cooperative computing in the parallel model of computation. In particular, it studies the *Do-All* problem in the shared-memory model, where it is referred to as the *Write-All* problem. The current monograph deals with the message-passing models of computation and considers broader adversarial settings inherent to these distributed models. The two monographs follow similar presentation philosophies and it is reasonable to consider them as complementary volumes. The current volume includes in Chapter 11 several recent results on the *Write-All* problem that appeared since the publication of the first monograph [67].

The book by Lynch [79] provides a wealth of information on distributed computing issues, such as computational models, algorithms, fault-tolerance, lower bounds and impossibility results. This include the consensus problem, which is related to *Do-All*, and we discuss this relation in Chapter 11. Additionally, information on the Input/Output Automata used in our Chapter 8 can be found there. The book by Attiya and Welch [6] is another excellent source of information on distributed computing issues, including cooperation. The book of Guerraoui and Rodrigues [52] presents numerous important abstractions for reliable distributed computing and includes detailed examples of how these abstractions can be implemented and used in practice.

Acknowledgements

Our research on robust distributed cooperation continues to be inspired by the earlier work on fault-tolerant parallel computing of the late Paris Christos Kanellakis (1953-1995). Paris is survived by his parents, General Eleftherios and Roula Kanellakis, who have enthusiastically encouraged us to continue his work through the long years following the tragic death of Paris, his wife Maria-Teresa, and their children Alexandra and Stephanos. We warmly thank General and Mrs. Kanellakis for their inspiration and support.

The material presented in this monograph includes results obtained by the authors in collaboration with Bogdan Chlebus, Roberto De Prisco, Antonio Fernandez, Dariusz Kowalski, Greg Malewicz and Alexander Russell. We thank them for the wonderful and fruitful collaborations—without their contributions this monograph would not exist.

Our work on robust distributed cooperation also benefited from prior collaboration with several colleagues, and we gratefully acknowledge the contributions of Jonathan Buss, Shlomi Dolev, Leszek Gasieniec, Dimitrios Michailidis, Prabhakar Ragde, and Roberto Segala.

Special thanks are due to Nancy Lynch for reviewing an earlier version of this work. Her insight and valuable feedback are greatly appreciated.

In undertaking the research that ultimately resulted in this monograph, we were motivated by the work of other researchers who have also contributed to the field of fault-tolerant cooperative computing. We will be remiss without mentioning the names of Richard Anderson, Richard Cole, Cynthia Dwork, Zvi Galil, Phillip Gibbons, Joe Halpern, Maurice Herlihy, Zvi Kedem, Andrzej Lingas, Chip Martel, Keith Marzullo, Alan Mayer, Naomi Nishimura, Krishna Palem, Arvin Park, Michael Rabin, Arvind Raghunathan, Nir Shavit, Paul Spirakis, Ramesh Subramonian, Orli Waarts, Heather Woll, Moti Yung, Ofer Zajicek, and Asaph Zemach.

This work was in part supported by the National Science Foundation (NSF) Grants 9984778, 9988304, 0121277, 0311368, and by the NSF-NATO Award 0209588. The work of the first author has also been partially supported by research funds from the University of Cyprus.

We thank our Springer editor, Susan Lagerstrom-Fife, for her encouragement and support, and Sharon Palleschi, editorial assistant at Springer, for her valuable assistance during the preparation of this monograph.

Our warmest thanks go to Michel Raynal for writing the foreword of this monograph; thank you for this honor, Michel.

Finally, we would like to thank our families.

CG: I thank my wife, Agni, for the emotional support she has given me and the patience she has shown during the endless nights I spent working on this monograph, while she took care of our angel, 21 months old son Yiorgo. You both bring joy and meaning to my life and you are my source of strength and inspiration. Agni, you are the *love* of my life. Yiorgo, you *are* my life.

AAS: My wife Sana gave me more affection, care, and happiness, than I could have dreamed of. She spent many lonely nights being an epitome of patience, while sustaining herself only by my assurances that I'll belong to her yet again. Thank you, my love. I thank my children, Ginger and Ted, for being there for me when I needed you most. This time around you are grown-ups and I am proud of you. I thank my step-son Arnold for his encouragement and interest. I never thought that daily questions from a freshman "Are you done with the book yet?" would do so much to energize the work of this professor. I am glad to report: "We are done."

Nicosia, Cyprus and Storrs, CT, USA
September 2007

Chryssis Georgiou
Alexander A. Shvartsman