2001-672

Amnon Barak Shai Guday Richard G. Wheeler

The MOSIX Distributed Operating System

Load Balancing for UNIX

Springer-Verlag

Berlin Heidelberg New York London Paris Tokyo Hong Kong Barcelona Budapest Series Editors

Gerhard Goos Universität Karlsruhe Postfach 69 80 Vincenz-Priessnitz-Straße 1 W-7500 Karlsruhe, FRO Juris Hartmanis Cornell University Department of Computer Science 4130 Upson Hall Ithaca, NY 14853, USA

Authors

Amnon Barak Shai Guday Richard G. Wheeler Institute of Computer Science, The Hebrew University of Jerusalem 91904 Jerusalem, Israel

CR Subject Classification (1991): D.4, C.2.4



ISBN 3-540-56663-5 Springer-Verlag Berlin Heidelberg New York ISBN 0-387-56663-5 Springer-Verlag New York Berlin Heidelberg

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

© Springer-Verlag Berlin Heidelberg 1993 Printed in Germany

Typesetting: Camera ready by authors/editors 45/3140-543210 - Printed on acid-free paper

Dedicated to our families.

Preface

This book describes the design and internals of the MOSIX distributed operating system. MOSIX, an acronym for Multicomputer Operating System for UNIX, integrates a cluster of loosely connected computers into a virtual single-machine UNIX environment. The main property of MOSIX is the *high degree of integration* among the computers, which may include personal workstations, shared memory and non-shared memory multiprocessors, connected by fast communication links. This integration includes network transparency, cooperation between the computers to provide services across machine boundaries, support of dynamic configuration, and system-initiated load balancing by process migration. Another property of MOSIX is the ability to scale up the system configuration to encompass a large number of computers. This is accomplished by using probabilistic algorithms that allow each computer to maintain only partial knowledge about the state of the global system, regardless of the number of computers.

The development of MOSIX was begun in 1981 for a cluster of PDP-11 computers. It was based on UNIX Version 7. Since then, four additional versions of MOSIX have been developed, each version based on the most recent version of UNIX that was available at the time. The latest version is operational on a cluster of workstations, where each workstation is itself a multiprocessor.

We describe MOSIX as it is, rather than as it was meant to be, or as it would be implemented today. The text summarizes some relevant parts of UNIX that provide a basis for understanding MOSIX. Readers interested in detailed descriptions of UNIX should refer to any of the excellent texts referred to in this book. The material presented is intended primarily for readers who are interested in distributed and multiprocessor systems. The reader is assumed to have some knowledge in programming and operating systems, preferably UNIX. Readers without this background will still benefit from the techniques and algorithms discussed.

The book consists of eleven chapters. Chapter 1 gives a brief introduction to MOSIX, placing it in the context of other multicomputer systems. Chapter 2 presents an overview of MOSIX, its characteristics, kernel architecture, and developmental stages. Chapter 3 describes the design of the traditional, non-distributed UNIX file system, followed by Chapter 4, which describes several distributed UNIX file systems and the MOSIX file system in particular. Chapter 5 describes the non-distributed UNIX process structure and the kernel mechanisms that interact with the processes. Chapter 6 presents the MOSIX process, the internal mechanisms that allow process migration in MOSIX, and the details of the distributed interprocess communication mechanisms. The underlying links between the host-specific part of the kernel and the user view of the kernel are presented in Chapter 7.

Chapter 8 describes the load balancing mechanism built into the MOSIX kernel. Scaling considerations, including examples of probabilistic algorithms that are used in MOSIX, are presented in Chapter 9. Chapter 10 presents the performance of the main communication mechanisms of the MOSIX kernel. Chapter 11 discusses distributed applications. It presents a brief description of a language extension for writing distributed programs and gives the performance results of several examples. It concludes with a description of a monitoring facility for distributed applications that is supported by extensions to the MOSIX kernel.

The development of MOSIX has been a cooperative effort of many individuals, but some call for particular mention here. First, we must acknowledge the significant contributions of Amnon Shiloh to the design and implementation of all the MOSIX versions, and for reviewing drafts of the manuscript. Ami Litman made valuable contributions to the design and development of the first version of MOSIX. Special thanks to Danny Braniss for his help and to Avi Barel for the implementation of NSMOS. Thanks are also due to Robert Hofner, Roy Laor, Jonathan Masel, On G. Paradise, Gil Shwed, Yuval Yarom and all of the other participants on the MOSIX team for their contributions. We would like to thank National Semiconductor Corporation (Israel) for their support and equipment contributions, and Jennifer Steiner for editing the manuscript.

The research that led to the development of MOSIX was supported in part by the U.S. Air Force, Office of Scientific Research, sponsored by the HQ Rome Air Development Center and the European Office of Aerospace Research and Development, the Israel Ministry of Science and Technology, the Israel National Council for Higher Education, the Israel Ministry of Defense, National Semiconductor Corporation, and the Israel Academy of Science and Humanities.

The following terms are trademarks: Ethernet (Xerox Corporation), M68000 (Motorola Semiconductor Corporation); NFS (Sun Microsystems, Inc.); PDP-11, Q-Bus, VAX, and VMS (Digital Equipment Corporation); ProNET, ProNET-10, ProNET-80 (Proteon, Inc.); VME532, VR32, NS32000, NS32332, and NS32532 (National Semiconductor Corporation); and UNIX (UNIX System Laboratories).

A. Barak, S. Guday, R.G. Wheeler Jerusalem, 1993

Contents

1	Introduction	1
2	Overview of MOSIX	5
	2.1 The Characteristics of MOSIX	5
	2.2 The Architecture of the MOSIX Kernel	9
	2.3 The History of the MOSIX Project	2
	2.4 Summary 1	.7
3	The UNIX File System 1	9
	3.1 The Namespace	20
	3.2 The Traditional File System	21
	3.3 UNIX Buffer Caching	26
	3.4 UNIX File System Calls	32
	3.5 Summary	36
4	Distributed UNIX File Systems 3	17
	4.1 Extending the Traditional Namespace	37
	4.2 Classifying Distributed File Systems	12
	4.3 MOSIX File System Implementation	13
	4.4 MOSIX File System Calls	19
	4.5 Summary	76
5	The UNIX Process 7	'7
	5.1 Organization of the System Memory	78
	5.2 Organization of the Process	31
	5.3 Process Context	34
	5.4 Process States	37
	5.5 Scheduling Processes	38
	5.6 Process System Calls	91
	5.7 Summary	38
5	4.5 Summary 7 The UNIX Process 7 5.1 Organization of the System Memory 7 5.2 Organization of the Process 8 5.3 Process Context 8 5.4 Process States 8 5.5 Scheduling Processes 8 5.6 Process System Calls 9 5.7 Summary 9	73333

6.1 Remote Paging 99 6.2 MOSIX Process Structure 100 6.3 MOSIX Process System Calls 103 6.4 Process Migration 109 6.5 Interprocess Communication 113 6.6 Summary 114 7 The Interface Layer 115 7.1 The Interface Layer 122 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 136 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing Policy 143 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Considerations 4lgorithms 160 8.9 Summary 167 9 Scaling Considerations in MOSIX 172 9.1 <t< th=""><th>6</th><th>\mathbf{The}</th><th>MOSIX Process</th><th>99</th></t<>	6	\mathbf{The}	MOSIX Process	99		
6.2 MOSIX Process Structure 100 6.3 MOSIX Process System Calls 103 6.4 Process Migration 103 6.5 Interprocess Communication 113 6.6 Summary 114 7 The MOSIX Linker 115 7.1 The Interface Layer 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.2 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance		6.1	Remote Paging	- 99		
6.3 MOSIX Process System Calls 103 6.4 Process Migration 109 6.5 Interprocess Communication 113 6.6 Summary 114 7 The MOSIX Linker 115 7.1 The Interface Layer 114 7 The MOSIX Linker 115 7.1 The Interface Layer 114 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 136 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing Policy 143 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 152 8.8 The Migration Consideration Algorithms 152 8.8 The Migration Consideration Algorithms 160 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.2 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10.1 Scall Performance 180 10.2 Funnel & Process M		6.2	MOSIX Process Structure	100		
6.4 Process Migration 109 6.5 Interprocess Communication 113 6.6 Summary 114 7 The MOSIX Linker 115 7.1 The Interface Layer 117 7.2 The Transport Layer 122 7.4 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 152 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 169 9.2 Scaling Considerations in MOSIX 172 9.3 Probabilistic		6.3	MOSIX Process System Calls	103		
6.5 Interprocess Communication 113 6.6 Summary 114 7 The MOSIX Linker 115 7.1 The Interface Layer 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 140 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Information Algorithms 160 8.9 Summary 167 9 Scaling Considerations MOSIX 172 9.3 Probabilistic Algorithms 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 169 9.1 Summary		6.4	Process Migration	109		
6.6 Summary 114 7 The MOSIX Linker 115 7.1 The Interface Layer 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 147 8.7 The Information Consideration Algorithms 152 8.8 The Migration Considerations 160 8.9 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 182 10.1 </td <td></td> <td>6.5</td> <td>Interprocess Communication</td> <td>113</td>		6.5	Interprocess Communication	113		
7 The MOSIX Linker 115 7.1 The Interface Layer 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Laod Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.2 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 180 10.2 Funnel & Process Migration Performance 183 10.4 DAEMON T		6.6	Summary	114		
7 The MOSIX Linker 115 7.1 The Interface Layer. 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 140 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 180 10.2 Funnel & Process M						
7.1 The Interface Layer 117 7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 152 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 189 10.1 Scall Performance 180 10.2 Funnel & Process Migration Performance 182 10.3 Load Balancing Performance 183 10	7	The	MOSIX Linker	115		
7.2 The Transport Layer 124 7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 136 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing Policy 143 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 166 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 169 9.1 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 180 10.2 Funnel & Process Migration Performance 182		7.1	The Interface Layer	117		
7.3 The Network Layer 132 7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 139 8.4 Pre-emptive Load Balancing 140 8.4 Pre-emptive Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 179 10.1 Scall Performance 180 10.2 Funnel & Process Migration Performance 182 10.3 Load Balancing Performance 182 10.4<		7.2	The Transport Layer	124		
7.4 Summary 134 8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 139 8.4 Pre-emptive Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 160 8.9 Summary 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.2 Scaling Considerations in MOSIX 172 9.3 Probabilistic Algorithms 175 9.4 Summary 178 10 System Performance 179 10.1 Scall Performance 180 10.2 Funnel & Process Migration Performance 182 10.3 Load Balancing Performance 183 <tr< td=""><td></td><td>7.3</td><td>The Network Layer</td><td>132</td></tr<>		7.3	The Network Layer	132		
8 Load Balancing 135 8.1 Foundations of Load Balancing 136 8.2 Static Load Balancing 139 8.3 Dynamic Load Balancing 140 8.4 Pre-emptive Load Balancing 142 8.5 The MOSIX Load Balancing Policy 143 8.6 The Load Calculation Algorithms 147 8.7 The Information Dissemination Algorithms 152 8.8 The Migration Consideration Algorithms 152 8.8 The Migration Consideration Algorithms 167 9 Scaling Considerations 169 9.1 Principles of Scaling 169 9.1 Principles of Scaling 172 9.3 Probabilistic Algorithms 172 9.3 Probabilistic Algorithms 178 10 System Performance 180 10.2 Funnel & Process Migration Performance 182 10.3 Load Balancing Performance 183 10.4 DAEMON Toolkit Performance 185 10.5 Summary 189 11.1 Writing Distributed		7.4	Summary	134		
8.1Foundations of Load Balancing1368.2Static Load Balancing1398.3Dynamic Load Balancing1408.4Pre-emptive Load Balancing1428.5The MOSIX Load Balancing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18310.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217	8	Loa	d Balancing	135		
8.2Static Load Balancing1398.3Dynamic Load Balancing1408.4Pre-emptive Load Balancing1428.5The MOSIX Load Balancing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18311Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications19611.4Summary210Bibliography213Index217	Ť	8.1	Foundations of Load Balancing	136		
8.3Dynamic load Balancing1408.4Pre-emptive Load Balancing1428.5The MOSIX Load Balancing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18310.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.2Examples of Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		89	Static Load Balancing	130		
8.3Dynamic boar balancing1428.4Pre-emptive Load Balancing1428.5The MOSIX Load Balancing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications19611.2Examples of Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		8.2	Duramia Load Balanging	100		
8.4Treenprive Load Balancing Policy1428.5The MOSIX Load Balancing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		0.0 V A	Pro emptivo Load Balancing	140		
8.3The MOSIA Load Balatcing Policy1438.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scal Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18310.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		0.4 0 E	The MOCIV Load Delensing Delive	142		
8.6The Load Calculation Algorithms1478.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		8.0	The MOSIA Load Balancing Policy	143		
8.7The Information Dissemination Algorithms1528.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.2Examples of Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary213Bibliography213Index217		8.6	The Load Calculation Algorithms	147		
8.8The Migration Consideration Algorithms1608.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary213Bibliography213Index217		8.7	The Information Dissemination Algorithms	152		
8.9Summary1679Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary213Bibliography213Index217		8.8	The Migration Consideration Algorithms	160		
9Scaling Considerations1699.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		8.9	Summary	167		
9.1Principles of Scaling1699.2Scaling Considerations in MOSIX1729.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217	9	Scal	ing Considerations	169		
9.11.11.11.11.19.2Scaling Considerations in MOSIX1.729.3Probabilistic Algorithms1.759.4Summary1.759.4Summary1.7810System Performance1.7910.1Scall Performance1.8010.2Funnel & Process Migration Performance1.8210.3Load Balancing Performance1.8310.4DAEMON Toolkit Performance1.8510.5Summary1.8711Distributed Applications1.8911.1Writing Distributed Applications1.9611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217	•	91	Principles of Scaling	169		
9.3Probabilistic Algorithms1759.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18510.5Summary18711Distributed Applications18911.1Writing Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		9.2	Scaling Considerations in MOSIX	179		
9.4Summary17810System Performance17910.1Scall Performance18010.2Funnel & Process Migration Performance18210.3Load Balancing Performance18310.4DAEMON Toolkit Performance18310.5Summary18711Distributed Applications18911.1Writing Distributed Applications18911.2Examples of Distributed Applications19611.3Monitoring Distributed Applications20711.4Summary210Bibliography213Index217		0.3	Probabilistic Algorithms	175		
10 System Performance17910.1 Scall Performance18010.2 Funnel & Process Migration Performance18210.3 Load Balancing Performance18310.4 DAEMON Toolkit Performance18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		0.A		178		
10 System Performance17910.1 Scall Performance18010.2 Funnel & Process Migration Performance18210.3 Load Balancing Performance18310.4 DAEMON Toolkit Performance18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		<i>J</i> .7	Building (, , , , , , , , , , , , , , , , , ,	110		
10.1 Scall Performance18010.2 Funnel & Process Migration Performance18210.3 Load Balancing Performance18310.4 DAEMON Toolkit Performance18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217	10	Syst	tem Performance	179		
10.2 Funnel & Process Migration Performance18210.3 Load Balancing Performance18310.4 DAEMON Toolkit Performance18510.5 Summary18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		10.1	Scall Performance	180		
10.3 Load Balancing Performance18310.4 DAEMON Toolkit Performance18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		10.2	Funnel & Process Migration Performance	182		
10.4 DAEMON Toolkit Performance18510.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		10.3	Load Balancing Performance	183		
10.5 Summary18711 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		10.4	DAEMON Toolkit Performance	185		
11 Distributed Applications18911.1 Writing Distributed Applications18911.2 Examples of Distributed Applications19611.3 Monitoring Distributed Applications20711.4 Summary210Bibliography213Index217		10.5	Summary	187		
11.1 Writing Distributed Applications 189 11.2 Examples of Distributed Applications 196 11.3 Monitoring Distributed Applications 207 11.4 Summary 210 Bibliography 213 Index 217	11 Distributed Applications					
11.1 Writing Distributed Applications 169 11.2 Examples of Distributed Applications 196 11.3 Monitoring Distributed Applications 207 11.4 Summary 210 Bibliography 213 Index 217	τ.	11.1	Writing Distributed Applications	180		
11.2 Examples of Distributed Applications 136 11.3 Monitoring Distributed Applications 207 11.4 Summary 210 Bibliography 213 Index 217		11.1	Framples of Distributed Applications	106		
11.3 Monitoring Distributed Applications 207 11.4 Summary 210 Bibliography 213 Index 217		11.2	Monitoring Distributed Applications	100		
Bibliography 213 Index 217		11.0	Summer Distributed Applications	207		
Bibliography213Index217		11.4	Summary	210		
Index 217	Bi	Bibliography				
	In	Index				

Chapter 1

Introduction

Configurations of loosely-coupled multicomputers can be classified as either network systems or distributed systems. In a network system, each computer runs its own operating system. Each of these operating systems is augmented by communication facilities that permit interaction with the other systems in the network. In a network system, the user's environment is confined to one (local) computer, with the added ability to access objects that reside on other (remote) computers. This last ability, however, is quite limited, since all network commands must specify the location of each remote object. For example, users access remote objects (e.g., reading a file that resides on another computer) by explicitly using the network name of the computer that has the file. As a result, network operating systems have a limited degree of resource sharing and they are therefore commonly used to connect geographically dispersed and heterogeneous systems.

Distributed systems provide a higher degree of transparency and resource sharing than network systems. Distributed systems can be divided into two categories: user-level distributed systems and distributed operating systems. In a user-level distributed system, the support for distribution is provided in a layer of software on top of the (non-distributed) operating system. User-level distributed systems are intended for configurations of machines running different operating systems with the distributed software layer on top. An example of this type of distributed system is the Open Software Foundation's Distributed Computing Environment [17]. Distributed operating systems, on the other hand, implement support for distribution in the kernel, and are intended for configurations of machines running the same operating system. The MOSIX system described in this book is an example of a distributed operating system.

In a distributed operating system, one operating system is used by all the computers in the entire network, each computer running its own copy. Distributed operating systems are most commonly used in networks in which all of the computers are from the same manufacturer. Distributed operating systems can also be used in networks of computers from different manufacturers if all of the operating systems have the same functionality (i.e., the same user interface). The main goal of distributed operating systems is to provide resource sharing over a transparent network. In these systems, the user is provided with a single virtual machine, with transparent network communication, distribution of the workload, and automatic resource allocation.

The architecture of many distributed operating systems is based on the **client/server model**. Among these distributed operating systems, it is possible to distinguish between asymmetrical client/server systems and symmetrical client/server systems. In an asymmetrical client/server system, specific machines are assigned service functions (e.g., file servers or name servers), while other machines are used to execute user tasks. Examples of asymmetrical client/server distributed operating systems include the V-System [14] and Amoeba [29].

The main characteristic of a symmetrical client/server system is the decentralization of control; each machine is both a server and a client for all of the services. In symmetrical systems, each machine can function as an independent computer, with complete hardware and software facilities, while the entire network behaves like a single computer. The advantages of symmetrical systems over asymmetrical systems include improved cost/performance ratio, better resource utilization, increased availability and reliability, and the possibility to scale up the configuration to large numbers of computers.

The MOSIX system is a symmetrical distributed operating system that integrates a cluster of loosely connected, independent computers into a virtual single-machine UNIX environment. The hardware configuration for MOSIX consists of a cluster of computers, each with its own local memory, that are loosely connected by a local area communication network (LAN). In most configurations that have been developed to run MOSIX, each computer (node) is an independent uniprocessor UNIX system, with complete hardware and software facilities. In the latest configuration running MOSIX, each node is itself a multiprocessor. Each such node may contain up to eight independent processors that share I/O devices and communication controllers over a common bus. In this book, the terms "node" machine, and "workstation" are used to refer to an independent computer. The term "processor" refers to a single Processing Element (PE) within a multiprocessor workstation.

The main characteristics of MOSIX are:

Network transparency - the network is completely invisible to the naive user.

Autonomy each node is capable of operating as an independent system.

Cooperation the nodes work together to provide services across the network.

Decentralized control – each processor makes all of its control decisions independently.

Dynamic process migration - processes can be migrated among homogeneous processors.

З,

- Load balancing process migration allows near-optimal assignment of processes to processors.
- **Dynamic configuration** nodes may be added and removed with minimal side effects.
- Increased availability files and processes can be replicated on different nodes.
- **Performance** local and remote operations are highly efficient.
- **Reliability** a limited degree of reliability is provided through the isolation of faults.

Replicated kernel and resources - the system is replicated in each node.

Scalability \cdot the configuration may be scaled up to a large number of nodes.

Compatibility - MOSIX is compatible with AT&T UNIX System V.

More details about these characteristics are given in Chapter 2.

The most noticeable properties for executing distributed applications on MOSIX are its network transparency, the symmetry and flexibility of its configuration, and its dynamic process migration. The combined effect of these properties is that application programs are completely independent of the current state of the system configuration. Users do not need to change their applications due to node or communication links failures, nor be concerned about the load of the various processors. The system automatically attempts to optimize all resource allocation, including migrating I/O-bound jobs to the sites that master the devices they use and migrating heavily communicating jobs to nearby processors so that they benefit from fast services.

The MOSIX kernel is obtained by restructuring the UNIX kernel into machine-dependent and machine-independent parts (modules). The kernel is built as a structure of loosely-coupled modules, where the module interfaces are minimal and well defined [6]. The machine-dependent module provides sitedependent services, such as access to local disks. The machine-independent module provides network-wide services to the application level, such as transparent interprocessor communication. Each MOSIX kernel isolates the users' processes from the specific machines on which they execute, while at the same time providing these processes with the standard UNIX interface [2, 3]. This means that processes execute in a site-independent mode, which allows all system calls to be executed uniformly, regardless of the current location of the requesting process and the site that has the requested object.

The MOSIX kernel is designed to hide the internal network from the user and the application programmer [6]. In order to provide efficient networkwide services, MOSIX kernels interact with each other at the level of kernel remote procedure calls. The MOSIX kernel can be implemented on any reasonable hardware, but the participating processors must be homogeneous, to allow process-migration. This does not exclude MOSIX from being part of wider, heterogeneous networks, where process migration occurs among disjoined sets of homogeneous processors, improving resource sharing and performance.

Another major objective of the MOSIX kernel architecture is high performance. An important consideration in the design of this kernel was the avoidance of the high overhead caused by information hiding. Another architectural consideration was to enable comprehensive debugging of the system on a single machine. This means that, if necessary, any of the known debugging techniques can be applied to the MOSIX kernel. The details of the MOSIX kernel architecture are given in Chapter 2.

MOSIX belongs to a class of modular operating systems, designed as a set of system servers on top of a **microkernel** that provides low-level services (e.g., memory management and interprocess communication). In this type of operating system, the microkernel forms a standard base that can support higher level system-specific interfaces. The specific interface that is supported by MOSIX is UNIX. Examples of other microkernel-based operating systems include Amoeba [29], BirliX [20], Chorus [21], and Mach [46].

There are currently only a few distributed operating systems that support process migration. This is due to the difficulties of managing the migration itself, the need to change the operating system kernel architecture to support migration, the need to change the environment of the process (e.g., its dependence on local kernel tables), and the need to provide access to these resources from other sites. Examples of other operating systems that support process migration include Rhodos [49] and Sprite [15, 30]. For further reading about the design of distributed operating systems see [18]. A comprehensive survey of distributed systems is given in [12].

The motivation behind the MOSIX project is to research and develop operating systems for message passing-based (as opposed to shared memory) distributed systems. This research began in 1981, with the development of the original version of MOSIX for a cluster of PDP-11 computers. That version was compatible with UNIX Version 7 [47]. The second, M68000-based system was compatible with UNIX Version 7 with enhancements from BSD 4.1. Three additional versions have been developed for National Semiconductor's VR32, the VAX family, and National Semiconductor's VME532 multiprocessor architecture. These versions are compatible with UNIX System V Release 2 [2]. The development strategy for all these versions was to use as many ready-made hardware and software components as possible. This strategy enabled the development team to concentrate on such issues as the kernel architecture, probabilistic algorithms, performance, and debugging aids without spending too much effort on hardware design or the supporting software outside the operating system kernel.