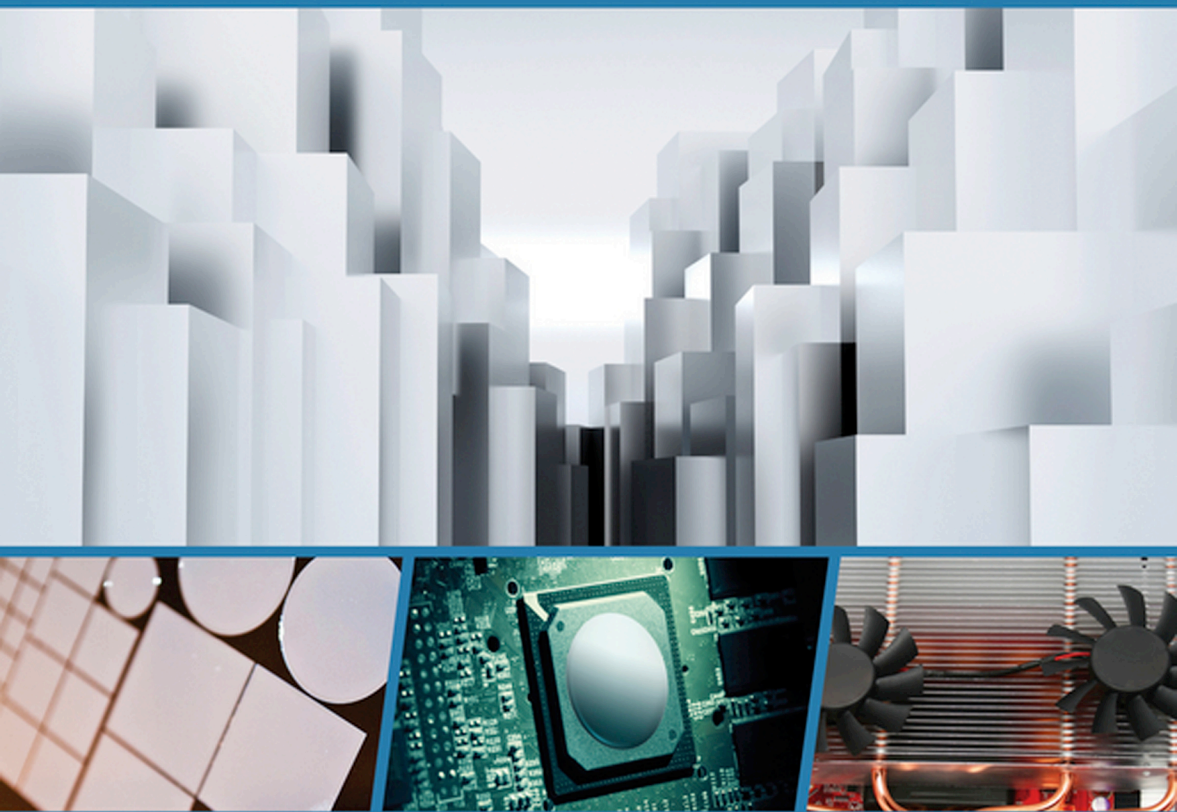


Wiley Series on Parallel and Distributed Computing

Albert Y. Zomaya, Series Editor



High-Performance Computing on Complex Environments

EDITED BY

Emmanuel Jeannot & Julius Žilinskas

WILEY

*High-Performance
Computing on Complex
Environments*

**WILEY SERIES ON PARALLEL
AND DISTRIBUTED COMPUTING**

Series Editor: Albert Y. Zomaya

A complete list of titles in this series appears at the end of this volume.

High-Performance Computing on Complex Environments

Emmanuel Jeannot

Inria

Julius Žilinskas

Vilnius University

WILEY

Copyright © 2014 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.

Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services please contact our Customer Care Department with the U.S. at 877-762-2974, outside the U.S. at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print, however, may not be available in electronic format.

Library of Congress Cataloging in Publication Data:

Jeannot, Emmanuel.

High performance computing on complex environments / Emmanuel Jeannot, Julius Žilinskas.
pages cm

Includes bibliographical references and index.

ISBN 978-1-118-71205-4 (cloth)

1. High performance computing. I. Žilinskas, J. (Julius), 1973- II. Title.

QA76.88.J43 2014

004.1'1-dc23

2013048363

High-Performance Computing on Complex Environments / Emmanuel Jeannot and Julius Žilinskas
Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

To our colleague Mark Baker

Contents

Contributors	xxiii
---------------------	--------------

Preface	xxvii
----------------	--------------

PART I INTRODUCTION	1
--------------------------------	----------

1. Summary of the Open European Network for High-Performance Computing in Complex Environments	3
---	----------

Emmanuel Jeannot and Julius Žilinskas

- 1.1 Introduction and Vision / 4
- 1.2 Scientific Organization / 6
 - 1.2.1 Scientific Focus / 6
 - 1.2.2 Working Groups / 6
- 1.3 Activities of the Project / 6
 - 1.3.1 Spring Schools / 6
 - 1.3.2 International Workshops / 7
 - 1.3.3 Working Groups Meetings / 7
 - 1.3.4 Management Committee Meetings / 7
 - 1.3.5 Short-Term Scientific Missions / 7
- 1.4 Main Outcomes of the Action / 7
- 1.5 Contents of the Book / 8
- Acknowledgment / 10

PART II	NUMERICAL ANALYSIS FOR HETEROGENEOUS AND MULTICORE SYSTEMS	11
2.	On the Impact of the Heterogeneous Multicore and Many-Core Platforms on Iterative Solution Methods and Preconditioning Techniques	13
	<i>Dimitar Lukarski and Maya Neytcheva</i>	
2.1	Introduction / 14	
2.2	General Description of Iterative Methods and Preconditioning / 16	
2.2.1	Basic Iterative Methods / 16	
2.2.2	Projection Methods: CG and GMRES / 18	
2.3	Preconditioning Techniques / 20	
2.4	Defect-Correction Technique / 21	
2.5	Multigrid Method / 22	
2.6	Parallelization of Iterative Methods / 22	
2.7	Heterogeneous Systems / 23	
2.7.1	Heterogeneous Computing / 24	
2.7.2	Algorithm Characteristics and Resource Utilization / 25	
2.7.3	Exposing Parallelism / 26	
2.7.4	Heterogeneity in Matrix Computation / 26	
2.7.5	Setup of Heterogeneous Iterative Solvers / 27	
2.8	Maintenance and Portability / 29	
2.9	Conclusion / 30	
	Acknowledgments / 31	
	References / 31	
3.	Efficient Numerical Solution of 2D Diffusion Equation on Multicore Computers	33
	<i>Matjaž Depolli, Gregor Kosec, and Roman Trobec</i>	
3.1	Introduction / 34	
3.2	Test Case / 35	
3.2.1	Governing Equations / 35	
3.2.2	Solution Procedure / 36	
3.3	Parallel Implementation / 39	
3.3.1	Intel PCM Library / 39	
3.3.2	OpenMP / 40	

- 3.4 Results / 41
 - 3.4.1 Results of Numerical Integration / 41
 - 3.4.2 Parallel Efficiency / 42
- 3.5 Discussion / 45
- 3.6 Conclusion / 47
 - Acknowledgment / 47
 - References / 47

4. Parallel Algorithms for Parabolic Problems on Graphs in Neuroscience **51**

Natalija Tumanova and Raimondas Čiegis

- 4.1 Introduction / 51
- 4.2 Formulation of the Discrete Model / 53
 - 4.2.1 The θ -Implicit Discrete Scheme / 55
 - 4.2.2 The Predictor–Corrector Algorithm I / 57
 - 4.2.3 The Predictor–Corrector Algorithm II / 58
- 4.3 Parallel Algorithms / 59
 - 4.3.1 Parallel θ -Implicit Algorithm / 59
 - 4.3.2 Parallel Predictor–Corrector Algorithm I / 62
 - 4.3.3 Parallel Predictor–Corrector Algorithm II / 63
- 4.4 Computational Results / 63
 - 4.4.1 Experimental Comparison of Predictor–Corrector Algorithms / 66
 - 4.4.2 Numerical Experiment of Neuron Excitation / 68
- 4.5 Conclusions / 69
 - Acknowledgments / 70
 - References / 70

PART III COMMUNICATION AND STORAGE CONSIDERATIONS IN HIGH-PERFORMANCE COMPUTING **73**

5. An Overview of Topology Mapping Algorithms and Techniques in High-Performance Computing **75**

Torsten Hoefler, Emmanuel Jeannot, and Guillaume Mercier

- 5.1 Introduction / 76
- 5.2 General Overview / 76
 - 5.2.1 A Key to Scalability: Data Locality / 77

5.2.2	Data Locality Management in Parallel Programming Models / 77
5.2.3	Virtual Topology: Definition and Characteristics / 78
5.2.4	Understanding the Hardware / 79
5.3	Formalization of the Problem / 79
5.4	Algorithmic Strategies for Topology Mapping / 81
5.4.1	Greedy Algorithm Variants / 81
5.4.2	Graph Partitioning / 82
5.4.3	Schemes Based on Graph Similarity / 82
5.4.4	Schemes Based on Subgraph Isomorphism / 82
5.5	Mapping Enforcement Techniques / 82
5.5.1	Resource Binding / 83
5.5.2	Rank Reordering / 83
5.5.3	Other Techniques / 84
5.6	Survey of Solutions / 85
5.6.1	Algorithmic Solutions / 85
5.6.2	Existing Implementations / 85
5.7	Conclusion and Open Problems / 89
	Acknowledgment / 90
	References / 90

6. Optimization of Collective Communication for Heterogeneous HPC Platforms

95

Kiril Dichev and Alexey Lastovetsky

6.1	Introduction / 95
6.2	Overview of Optimized Collectives and Topology-Aware Collectives / 97
6.3	Optimizations of Collectives on Homogeneous Clusters / 98
6.4	Heterogeneous Networks / 99
6.4.1	Comparison to Homogeneous Clusters / 99
6.5	Topology- and Performance-Aware Collectives / 100
6.6	Topology as Input / 101
6.7	Performance as Input / 102
6.7.1	Homogeneous Performance Models / 103
6.7.2	Heterogeneous Performance Models / 105

6.7.3	Estimation of Parameters of Heterogeneous Performance Models / 106	
6.7.4	Other Performance Models / 106	
6.8	Non-MPI Collective Algorithms for Heterogeneous Networks / 106	
6.8.1	Optimal Solutions with Multiple Spanning Trees / 107	
6.8.2	Adaptive Algorithms for Efficient Large-Message Transfer / 107	
6.8.3	Network Models Inspired by BitTorrent / 108	
6.9	Conclusion / 111	
	Acknowledgments / 111	
	References / 111	
7.	Effective Data Access Patterns on Massively Parallel Processors	115
	<i>Gabriele Capannini, Ranieri Baraglia, Fabrizio Silvestri, and Franco Maria Nardini</i>	
7.1	Introduction / 115	
7.2	Architectural Details / 116	
7.3	K-Model / 117	
7.3.1	The Architecture / 117	
7.3.2	Cost and Complexity Evaluation / 118	
7.3.3	Efficiency Evaluation / 119	
7.4	Parallel Prefix Sum / 120	
7.4.1	Experiments / 125	
7.5	Bitonic Sorting Networks / 126	
7.5.1	Experiments / 131	
7.6	Final Remarks / 132	
	Acknowledgments / 133	
	References / 133	
8.	Scalable Storage I/O Software for Blue Gene Architectures	135
	<i>Florin Isaila, Javier Garcia, and Jesús Carretero</i>	
8.1	Introduction / 135	
8.2	Blue Gene System Overview / 136	
8.2.1	Blue Gene Architecture / 136	
8.2.2	Operating System Architecture / 136	

8.3	Design and Implementation / 138
8.3.1	The Client Module / 139
8.3.2	The I/O Module / 141
8.4	Conclusions and Future Work / 142
	Acknowledgments / 142
	References / 142

PART IV	EFFICIENT EXPLOITATION OF HETEROGENEOUS ARCHITECTURES	145
----------------	--	------------

9.	Fair Resource Sharing for Dynamic Scheduling of Workflows on Heterogeneous Systems	147
-----------	---	------------

Hamid Arabnejad, Jorge G. Barbosa, and Frédéric Suter

9.1	Introduction / 148
9.1.1	Application Model / 148
9.1.2	System Model / 151
9.1.3	Performance Metrics / 152
9.2	Concurrent Workflow Scheduling / 153
9.2.1	Offline Scheduling of Concurrent Workflows / 154
9.2.2	Online Scheduling of Concurrent Workflows / 155
9.3	Experimental Results and Discussion / 160
9.3.1	DAG Structure / 160
9.3.2	Simulated Platforms / 160
9.3.3	Results and Discussion / 162
9.4	Conclusions / 165
	Acknowledgments / 166
	References / 166

10.	Systematic Mapping of Reed–Solomon Erasure Codes on Heterogeneous Multicore Architectures	169
------------	--	------------

Roman Wyrzykowski, Marcin Wozniak, and Lukasz Kuczynski

10.1	Introduction / 169
10.2	Related Works / 171
10.3	Reed–Solomon Codes and Linear Algebra Algorithms / 172
10.4	Mapping Reed–Solomon Codes on Cell/B.E. Architecture / 173
10.4.1	Cell/B.E. Architecture / 173

10.4.2	Basic Assumptions for Mapping /	174
10.4.3	Vectorization Algorithm and Increasing its Efficiency /	175
10.4.4	Performance Results /	177
10.5	Mapping Reed–Solomon Codes on Multicore GPU Architectures /	178
10.5.1	Parallelization of Reed–Solomon Codes on GPU Architectures /	178
10.5.2	Organization of GPU Threads /	180
10.6	Methods of Increasing the Algorithm Performance on GPUs /	181
10.6.1	Basic Modifications /	181
10.6.2	Stream Processing /	182
10.6.3	Using Shared Memory /	184
10.7	GPU Performance Evaluation /	185
10.7.1	Experimental Results /	185
10.7.2	Performance Analysis using the Roofline Model /	187
10.8	Conclusions and Future Works /	190
	Acknowledgments /	191
	References /	191

11. Heterogeneous Parallel Computing Platforms and Tools for Compute-Intensive Algorithms: A Case Study **193**

Daniele D’Agostino, Andrea Clematis, and Emanuele Danovaro

11.1	Introduction /	194
11.2	A Low-Cost Heterogeneous Computing Environment /	196
11.2.1	Adopted Computing Environment /	199
11.3	First Case Study: The N -Body Problem /	200
11.3.1	The Sequential N -Body Algorithm /	201
11.3.2	The Parallel N -Body Algorithm for Multicore Architectures /	203
11.3.3	The Parallel N -Body Algorithm for CUDA Architectures /	204
11.4	Second Case Study: The Convolution Algorithm /	206
11.4.1	The Sequential Convolver Algorithm /	206
11.4.2	The Parallel Convolver Algorithm for Multicore Architectures /	207
11.4.3	The Parallel Convolver Algorithm for GPU Architectures /	208

11.5 Conclusions / 211

Acknowledgments / 212

References / 212

12. Efficient Application of Hybrid Parallelism in Electromagnetism Problems **215**

Alejandro Álvarez-Melcón, Fernando D. Quesada, Domingo Giménez, Carlos Pérez-Alcaraz, José-Ginés Picón, and Tomás Ramírez

12.1 Introduction / 215

12.2 Computation of Green's functions in Hybrid Systems / 216

12.2.1 Computation in a Heterogeneous Cluster / 217

12.2.2 Experiments / 218

12.3 Parallelization in Numa Systems of a Volume Integral Equation Technique / 222

12.3.1 Experiments / 222

12.4 Autotuning Parallel Codes / 226

12.4.1 Empirical Autotuning / 227

12.4.2 Modeling the Linear Algebra Routines / 229

12.5 Conclusions and Future Research / 230

Acknowledgments / 231

References / 232

PART V CPU + GPU COPROCESSING **235**

13. Design and Optimization of Scientific Applications for Highly Heterogeneous and Hierarchical HPC Platforms Using Functional Computation Performance Models **237**

David Clarke, Aleksandar Ilic, Alexey Lastovetsky, Vladimir Rychkov, Leonel Sousa, and Ziming Zhong

13.1 Introduction / 238

13.2 Related Work / 241

13.3 Data Partitioning Based on Functional Performance Model / 243

13.4 Example Application: Heterogeneous Parallel Matrix Multiplication / 245

13.5 Performance Measurement on CPUs/GPUs System / 247

13.6	Functional Performance Models of Multiple Cores and GPUs /	248
13.7	FPM-Based Data Partitioning on CPUs/GPUs System /	250
13.8	Efficient Building of Functional Performance Models /	251
13.9	FPM-Based Data Partitioning on Hierarchical Platforms /	253
13.10	Conclusion /	257
	Acknowledgments /	259
	References /	259
14.	Efficient Multilevel Load Balancing on Heterogeneous CPU + GPU Systems	261
	<i>Aleksandar Ilic and Leonel Sousa</i>	
14.1	Introduction: Heterogeneous CPU + GPU Systems /	262
14.1.1	Open Problems and Specific Contributions /	263
14.2	Background and Related Work /	265
14.2.1	Divisible Load Scheduling in Distributed CPU-Only Systems /	265
14.2.2	Scheduling in Multicore CPU and Multi-GPU Environments /	268
14.3	Load Balancing Algorithms for Heterogeneous CPU + GPU Systems /	269
14.3.1	Multilevel Simultaneous Load Balancing Algorithm /	270
14.3.2	Algorithm for Multi-Installment Processing with Multidistributions /	273
14.4	Experimental Results /	275
14.4.1	MSLBA Evaluation: Dense Matrix Multiplication Case Study /	275
14.4.2	AMPMD Evaluation: 2D FFT Case Study /	277
14.5	Conclusions /	279
	Acknowledgments /	280
	References /	280
15.	The All-Pair Shortest-Path Problem in Shared-Memory Heterogeneous Systems	283
	<i>Hector Ortega-Arranz, Yuri Torres, Diego R. Llanos, and Arturo Gonzalez-Escribano</i>	
15.1	Introduction /	283

15.2	Algorithmic Overview / 285
15.2.1	Graph Theory Notation / 285
15.2.2	Dijkstra's Algorithm / 286
15.2.3	Parallel Version of Dijkstra's Algorithm / 287
15.3	CUDA Overview / 287
15.4	Heterogeneous Systems and Load Balancing / 288
15.5	Parallel Solutions to The APSP / 289
15.5.1	GPU Implementation / 289
15.5.2	Heterogeneous Implementation / 290
15.6	Experimental Setup / 291
15.6.1	Methodology / 291
15.6.2	Target Architectures / 292
15.6.3	Input Set Characteristics / 292
15.6.4	Load-Balancing Techniques Evaluated / 292
15.7	Experimental Results / 293
15.7.1	Complete APSP / 293
15.7.2	512-Source-Node-to-All Shortest Path / 295
15.7.3	Experimental Conclusions / 296
15.8	Conclusions / 297
	Acknowledgments / 297
	References / 297

PART VI EFFICIENT EXPLOITATION OF DISTRIBUTED SYSTEMS

301

16. Resource Management for HPC on the Cloud

303

Marc E. Frincu and Dana Petcu

16.1	Introduction / 303
16.2	On the Type of Applications for HPC and HPC2 / 305
16.3	HPC on the Cloud / 306
16.3.1	General PaaS Solutions / 306
16.3.2	On-Demand Platforms for HPC / 310
16.4	Scheduling Algorithms for HPC2 / 311
16.5	Toward an Autonomous Scheduling Framework / 312
16.5.1	Autonomous Framework for RMS / 313

16.5.2	Self-Management / 315	
16.5.3	Use Cases / 317	
16.6	Conclusions / 319	
	Acknowledgment / 320	
	References / 320	
17.	Resource Discovery in Large-Scale Grid Systems	323
	<i>Konstantinos Karaogloulou and Helen Karatza</i>	
17.1	Introduction and Background / 323	
17.1.1	Introduction / 323	
17.1.2	Resource Discovery in Grids / 324	
17.1.3	Background / 325	
17.2	The Semantic Communities Approach / 325	
17.2.1	Grid Resource Discovery Using Semantic Communities / 325	
17.2.2	Grid Resource Discovery Based on Semantically Linked Virtual Organizations / 327	
17.3	The P2P Approach / 329	
17.3.1	On Fully Decentralized Resource Discovery in Grid Environments Using a P2P Architecture / 329	
17.3.2	P2P Protocols for Resource Discovery in the Grid / 330	
17.4	The Grid-Routing Transferring Approach / 333	
17.4.1	Resource Discovery Based on Matchmaking Routers / 333	
17.4.2	Acquiring Knowledge in a Large-Scale Grid System / 335	
17.5	Conclusions / 337	
	Acknowledgment / 338	
	References / 338	
PART VII	ENERGY AWARENESS IN HIGH-PERFORMANCE COMPUTING	341
18.	Energy-Aware Approaches for HPC Systems	343
	<i>Robert Basmadjian, Georges Da Costa, Ghislain Landry Tsafack Chetsa, Laurent Lefevre, Ariel Oleksiak, and Jean-Marc Pierson</i>	
18.1	Introduction / 344	
18.2	Power Consumption of Servers / 345	
18.2.1	Server Modeling / 346	

18.2.2	Power Prediction Models /	347
18.3	Classification and Energy Profiles of HPC Applications /	354
18.3.1	Phase Detection /	356
18.3.2	Phase Identification /	358
18.4	Policies and Leverages /	359
18.5	Conclusion /	360
	Acknowledgements /	361
	References /	361
19.	Strategies for Increased Energy Awareness in Cloud Federations	365
	<i>Gabor Kecskemeti, Attila Kertesz, Attila Cs. Marosi, and Zsolt Nemeth</i>	
19.1	Introduction /	365
19.2	Related Work /	367
19.3	Scenarios /	369
19.3.1	Increased Energy Awareness Across Multiple Data Centers within a Single Administrative Domain /	369
19.3.2	Energy Considerations in Commercial Cloud Federations /	372
19.3.3	Reduced Energy Footprint of Academic Cloud Federations /	374
19.4	Energy-Aware Cloud Federations /	374
19.4.1	Availability of Energy-Consumption-Related Information /	375
19.4.2	Service Call Scheduling at the Meta-Brokering Level of FCM /	376
19.4.3	Service Call Scheduling and VM Management at the Cloud-Brokering Level of FCM /	377
19.5	Conclusions /	379
	Acknowledgments /	380
	References /	380
20.	Enabling Network Security in HPC Systems Using Heterogeneous CMPs	383
	<i>Ozcan Ozturk and Suleyman Tosun</i>	
20.1	Introduction /	384
20.2	Related Work /	386

20.3	Overview of Our Approach / 387
20.3.1	Heterogeneous CMP Architecture / 387
20.3.2	Network Security Application Behavior / 388
20.3.3	High-Level View / 389
20.4	Heterogeneous CMP Design for Network Security Processors / 390
20.4.1	Task Assignment / 390
20.4.2	ILP Formulation / 391
20.4.3	Discussion / 393
20.5	Experimental Evaluation / 394
20.5.1	Setup / 394
20.5.2	Results / 395
20.6	Concluding Remarks / 397
	Acknowledgments / 397
	References / 397

PART VIII APPLICATIONS OF HETEROGENEOUS HIGH-PERFORMANCE COMPUTING 401

21. Toward a High-Performance Distributed CBIR System for Hyperspectral Remote Sensing Data: A Case Study in Jungle Computing 403

Timo van Kessel, Niels Drost, Jason Maassen, Henri E. Bal, Frank J. Seinstra, and Antonio J. Plaza

21.1	Introduction / 404
21.2	CBIR For Hyperspectral Imaging Data / 407
21.2.1	Spectral Unmixing / 407
21.2.2	Proposed CBIR System / 409
21.3	Jungle Computing / 410
21.3.1	Jungle Computing: Requirements / 411
21.4	IBIS and Constellation / 412
21.5	System Design and Implementation / 415
21.5.1	Endmember Extraction / 418
21.5.2	Query Execution / 418
21.5.3	Equi-Kernels / 419
21.5.4	Matchmaking / 420
21.6	Evaluation / 420
21.6.1	Performance Evaluation / 421

21.7 Conclusions / 426

Acknowledgments / 426

References / 426

22. Taking Advantage of Heterogeneous Platforms in Image and Video Processing 429

*Sidi A. Mahmoudi, Erencan Ozkan, Pierre Manneback,
and Suleyman Tosun*

22.1 Introduction / 430

22.2 Related Work / 431

22.2.1 Image Processing on GPU / 431

22.2.2 Video Processing on GPU / 432

22.2.3 Contribution / 433

22.3 Parallel Image Processing on GPU / 433

22.3.1 Development Scheme for Image Processing on GPU / 433

22.3.2 GPU Optimization / 434

22.3.3 GPU Implementation of Edge and Corner Detection / 434

22.3.4 Performance Analysis and Evaluation / 434

22.4 Image Processing on Heterogeneous Architectures / 437

22.4.1 Development Scheme for Multiple Image Processing / 437

22.4.2 Task Scheduling within Heterogeneous Architectures / 438

22.4.3 Optimization Within Heterogeneous Architectures / 438

22.5 Video Processing on GPU / 438

22.5.1 Development Scheme for Video Processing on GPU / 439

22.5.2 GPU Optimizations / 440

22.5.3 GPU Implementations / 440

22.5.4 GPU-Based Silhouette Extraction / 440

22.5.5 GPU-Based Optical Flow Estimation / 440

22.5.6 Result Analysis / 443

22.6 Experimental Results / 444

22.6.1 Heterogeneous Computing for Vertebra Segmentation / 444

22.6.2 GPU Computing for Motion Detection Using a Moving
Camera / 445

22.7 Conclusion / 447

Acknowledgment / 448

References / 448

23. Real-Time Tomographic Reconstruction Through CPU + GPU Coprocessing	451
<i>José Ignacio Agulleiro, Francisco Vazquez, Ester M. Garzon, and Jose J. Fernandez</i>	
23.1 Introduction /	452
23.2 Tomographic Reconstruction /	453
23.3 Optimization of Tomographic Reconstruction for CPUs and for GPUs /	455
23.4 Hybrid CPU + GPU Tomographic Reconstruction /	457
23.5 Results /	459
23.6 Discussion and Conclusion /	461
Acknowledgments /	463
References /	463
Index	467

Contributors

ALEJANDRO ÁLVAREZ-MELCÓN, Technical University of Cartagena, Cartagena, Spain
HAMID ARABNEJAD, Universidade do Porto, Porto, Portugal
HENRI E. BAL, VU University, Amsterdam, The Netherlands
RANIERI BARAGLIA, National Research Council of Italy, Pisa, Italy
JORGE G. BARBOSA, Universidade do Porto, Porto, Portugal
ROBERT BASMADJIAN, Passau University, Passau, Germany
GABRIELE CAPANNINI, D&IT Chalmers, Göteborg, Sweden
JESÚS CARRETERO, Universidad Carlos III of Madrid, Madrid, Spain
RAIMONDAS ČIEGIS, Vilnius Gediminas Technical University, Vilnius, Lithuania
DAVID CLARKE, University College Dublin, Dublin, Ireland
ANDREA CLEMATIS, IMATI CNR, Genoa, Italy
GEORGES DA COSTA, Toulouse University, Toulouse, France
DANIELE D'AGOSTINO, IMATI CNR, Genoa, Italy
EMANUELE DANOVARO, IMATI CNR, Genoa, Italy
MATJAŽ DEPOLLI, Jožef Stefan Institute, Ljubljana, Slovenia
KIRIL DICHEV, University College Dublin, Dublin, Ireland
NIELS DROST, Netherlands eScience Center, Amsterdam, The Netherlands
JOSE J. FERNANDEZ, National Centre for Biotechnology, National Research Council (CNB-CSIC), Madrid, Spain
MARC E. FRINCU, West University of Timisoara, Timisoara, Romania
JAVIER GARCIA, Universidad Carlos III of Madrid, Madrid, Spain
ESTER M. GARZON, University of Almería, Almería, Spain
DOMINGO GIMÉNEZ, University of Murcia, Murcia, Spain
ARTURO GONZALEZ-ESCRIBANO, Universidad de Valladolid, Valladolid, Spain
TORSTEN HOEFLER, ETH Zürich, Zürich, Switzerland

JOSÉ IGNACIO AGULLEIRO, University of Almería, Almería, Spain
ALEKSANDAR ILIC, Technical University of Lisbon, Lisbon, Portugal
FLORIN ISAILA, Universidad Carlos III of Madrid, Madrid, Spain
EMMANUEL JEANNOT, Inria Bordeaux Sud-Ouest, Talence, France
KONSTANTINOS KARAOGLANOGLU, Aristotle University of Thessaloniki, Thessaloniki, Greece
HELEN KARATZA, Aristotle University of Thessaloniki, Thessaloniki, Greece
GABOR KECSKEMETI, University of Innsbruck, Innsbruck, Austria
ATTILA KERTESZ, MTA SZTAKI Computer and Automation Research Institute, Budapest, Hungary
TIMO VAN KESSEL, VU University, Amsterdam, The Netherlands
GREGOR KOSEC, Jožef Stefan Institute, Ljubljana, Slovenia
LUKASZ KUCZYNSKI, Czestochowa University of Technology, Czestochowa, Poland
ALEXEY LASTOVETSKY, University College Dublin, Dublin, Ireland
LAURENT LEFEVRE, INRIA, LIP Laboratory, Ecole Normale Supérieure of Lyon, Lyon, France
DIEGO R. LLANOS, Universidad de Valladolid, Valladolid, Spain
DIMITAR LUKARSKI, Uppsala University, Uppsala, Sweden
JASON MAASSEN, Netherlands eScience Center, Amsterdam, The Netherlands
SIDI A. MAHMOUDI, University of Mons, Mons, Belgium
PIERRE MANNEBACK, University of Mons, Mons, Belgium
ATTILA Cs. MAROSI, MTA SZTAKI Computer and Automation Research Institute, Budapest, Hungary
GUILLAUME MERCIER, Bordeaux Polytechnic Institute, Talence, France; Inria Bordeaux Sud-Ouest, Talence, France
FRANCO MARIA NARDINI, National Research Council of Italy, Pisa, Italy
ZSOLT NEMETH, MTA SZTAKI Computer and Automation Research Institute, Budapest, Hungary
MAYA NEYTCHEVA, Uppsala University, Uppsala, Sweden
ARIEL OLEKSIK, Poznan Supercomputing and Networking Center, Poznan, Poland
HECTOR ORTEGA-ARRANZ, Universidad de Valladolid, Valladolid, Spain
ERENCAN OZKAN, Ankara University, Ankara, Turkey
OZCAN OZTURK, Bilkent University, Ankara, Turkey
CARLOS PÉREZ-ALCARAZ, University of Murcia, Murcia, Spain
DANA PETCU, West University of Timisoara, Timisoara, Romania
JOSÉ-GINÉS PICÓN, University of Murcia, Murcia, Spain
JEAN-MARC PIERSON, Toulouse University, Toulouse, France
FERNANDO D. QUESADA, Technical University of Cartagena, Cartagena, Spain
ANTONIO J. PLAZA, University of Extremadura, Caceres, Spain
TOMÁS RAMÍREZ, University of Murcia, Murcia, Spain
VLADIMIR RYCHKOV, University College Dublin, Dublin, Ireland
FRANK J. SEINSTR, Netherlands eScience Center, Amsterdam, The Netherlands
FABRIZIO SILVESTRI, National Research Council of Italy, Pisa, Italy
LEONEL SOUSA, Technical University of Lisbon, Lisbon, Portugal

FRÉDÉRIC SUTER, IN2P3 Computing Center, CNRS, IN2P3, Lyon-Villeurbanne, France

YURI TORRES, Universidad de Valladolid, Valladolid, Spain

SULEYMAN TOSUN, Ankara University, Ankara, Turkey

ROMAN TROBEC, Jožef Stefan Institute, Ljubljana, Slovenia

GHISLAIN LANDRY TSAFACK CHETSA, INRIA, LIP Laboratory, Ecole Normale Supérieure de Lyon, Lyon, France

NATALIJA TUMANOVA, Vilnius Gediminas Technical University, Vilnius, Lithuania

FRANCISCO VAZQUEZ, University of Almería, Almería, Spain

MARCIN WOZNIAK, Czestochowa University of Technology, Czestochowa, Poland

ROMAN WYRZYKOWSKI, Czestochowa University of Technology, Czestochowa, Poland

ZIMING ZHONG, University College Dublin, Dublin, Ireland

JULIUS ŽILINSKAS, Vilnius University, Vilnius, Lithuania

Preface

High-performance computing (HPC) is an important domain of the computer science field. For more than 30 years, it has allowed finding solutions to problems and enhanced progress in many scientific and industrial areas, such as climatology, biology, geology, and drug design, as well as automobile and aerospace engineering. However, new technologies such as multicore chips and accelerators have forced researchers in the field to rethink most of the advances in the domain, such as algorithms, runtime systems, language, software, and applications.

It is expected that a high-end supercomputer will be able to deliver several hundreds of petaflops (1 petaflop is 10^{15} floating-point operations per second) in 5 years from now. However, this will require mastering several challenges, such as energy efficiency, scalability, and heterogeneity.

Better and efficient parallel computers will enable solving problems at a scale and within a timeframe that has not been reached so far. These modern hierarchical and heterogeneous computing infrastructures are hard to program and use efficiently, particularly for extreme-scale computing. Consequently, none of the state-of-the-art solutions are able to efficiently use such environments. Providing tools for the whole software stack will allow programmers and scientists to efficiently write new program that will use most of the available power of such future complex machines.

COST Action IC0805 “Open European Network for High-Performance Computing on Complex Environments” (ComplexHPC) was devoted to heterogeneous and hierarchical systems for HPC, and is aimed at tackling the problem at every level (from cores to large-scale environments) and providing new integrated solutions for large-scale computing for future platforms. The duration of ComplexHPC Action was May 2009–June 2013. The goal of COST Action was to establish a European research network focused on high-performance heterogeneous computing to address the whole

range of challenges posed by these new platforms, including models, algorithms, programming tools, and applications. Indeed, some of the most active research groups in this area are in Europe. The network has contributed to exchanging information, identifying synergies, and pursuing common research activities, thereby reinforcing the strength of these groups and the leadership of Europe in this field. This book presents the results of COST Action. The chapters are written by expert participants of the Action.

This book is intended for scientists and researchers working in the field of HPC. It will provide advanced information for the readers already familiar with the basics of parallel and distributed computing. It may also be useful for PhD students and early stage researchers in computer science and engineering. It will also be of help to these young researchers to get a deep introduction to the related fields.

This book would not have been possible without the efforts of the contributors in preparing the respective chapters, and we would like to thank them for timely submissions and corrections. We would also like to thank Prof. Albert Zomaya for giving us the opportunity to publish this book in the “Wiley Series on Parallel and Distributed Computing.” We would also like to thank Simone Taylor, Director, Editorial Development, John Wiley & Sons, Inc., and the editorial team for their patience and guiding us through the publication of this book. We would also like to thank COST for the support that enabled the publication.

Delft, Netherlands
May, 2013

E. JEANNOT AND J. ŽILINSKAS



COST—the acronym for European Cooperation in Science and Technology—is the oldest and widest European intergovernmental network for cooperation in research. Established by the Ministerial Conference in November 1971, COST is presently used by the scientific communities of 36 European countries to cooperate in common research projects supported by national funds.

The funds provided by COST—less than 1% of the total value of the projects—support the COST cooperation networks (COST Actions) through which, with EUR 30 million per year, more than 30 000 European scientists are involved in research having a total value which exceeds EUR 2 billion per year. This is the financial worth of the European added value which COST achieves.

A “bottom up approach” (the initiative of launching a COST Action comes from the European scientists themselves), “à la carte participation” (only countries interested in the Action participate), “equality of access” (participation is open also to the scientific communities of countries not belonging to the European Union) and “flexible structure” (easy implementation and light management of the research initiatives) are the main characteristics of COST.

As precursor of advanced multidisciplinary research COST has a very important role for the realisation of the European Research Area (ERA) anticipating and complementing the activities of the Framework Programmes, constituting a “bridge” towards the scientific communities of emerging countries, increasing the mobility of researchers across Europe and fostering the establishment of “Networks of Excellence” in many key scientific domains such as: Biomedicine and Molecular Biosciences; Food and Agriculture; Forests, their Products and Services; Materials, Physical and Nanosciences; Chemistry and Molecular Sciences and Technologies; Earth System Science and Environmental Management; Information and Communication Technologies; Transport and Urban Development; Individuals, Societies, Cultures and Health. It covers basic and more applied research and also addresses issues of pre-normative nature or of societal importance.

Web: <http://www.cost.eu>

Neither the COST Office nor any person acting on its behalf is responsible for the use which might be made of the information contained in this publication. The COST Office is not responsible for the external websites referred to in this publication.

PART I

Introduction

1

Summary of the Open European Network for High-Performance Computing in Complex Environments

Emmanuel Jeannot

Inria Bordeaux Sud-Ouest, Talence, France

Julius Žilinskas

Vilnius University, Vilnius, Lithuania

In this chapter, we describe the COST Action IC0805 entitled “Open European Network for High-Performance Computing on Complex Environments.” This Action had representation from more than 20 countries and lasted from 2009 to 2013. We outline the scientific focus of this Action, its organization, and its main outcomes. The chapter concludes by presenting the structure of the book and its different chapters.

High-Performance Computing on Complex Environments, First Edition.

Edited by Emmanuel Jeannot and Julius Žilinskas.

© 2014 John Wiley & Sons, Inc. Published 2014 by John Wiley & Sons, Inc.

1.1 INTRODUCTION AND VISION

In recent years, the evolution and growth of the techniques and platforms commonly used for high-performance computing (HPC) in the context of different application domains has been truly astonishing. While parallel computing systems have now achieved certain maturity thanks to high-level libraries (such as ScaLAPACK) or runtime libraries (such as MPI), recent advances in these technologies pose several challenging research issues. Indeed, current HPC-oriented environments are extremely complex and very difficult to manage, particularly for extreme-scale application problems.

At the very low level, the latest generation CPUs are made of multicore processors that can be general-purpose or highly specialized in nature. On the other hand, several processors can be assembled into a so-called symmetrical multiprocessor (SMP) which can also have access to powerful specialized processors, such as graphics processing units (GPUs), that are now increasingly being used for programmable computing resulting from their advent in the video-game industry, which has significantly reduced their cost and availability. Modern HPC-oriented parallel computers are typically composed of several SMP nodes interconnected by a network. This kind of infrastructure is hierarchical and represents a first class of heterogeneous system in which the communication time between two processing units is different, depending on whether the units are on the same chip, on the same node, or not. Moreover, current hardware trends anticipate a further increase in the number of cores (in a hierarchical way) inside the chip, thus increasing the overall heterogeneity, even more toward building extreme-scale systems.

At a higher level, the emergence of heterogeneous computing now allows groups of users to benefit from networks of processors that are already available in their research laboratories. This is a second type of infrastructure where both the network and the processing units are heterogeneous in nature. Specifically, here the goal is to deal with networks that interconnect a (often high) number of heterogeneous computers that can significantly differ from one another in terms of their hardware and software architecture, including different types of CPUs operating at different clock speeds and under different design paradigms, and also different memory sizes, caching strategies, and operating systems.

At the high level, computers are increasingly interconnected together throughout wide area networks to form large-scale distributed systems with high computing capacity. Furthermore, computers located in different laboratories can collaborate in the solution of a common problem. Therefore, the current trends of HPC are clearly oriented toward extreme-scale, complex infrastructures with a great deal of intrinsic heterogeneity and many different hierarchical levels.

It is important to note that all the heterogeneity levels mentioned above are tightly linked. First of all, some of the nodes in computational distributed environments may be multicore SMP clusters. Second, multicore chips will soon be fully heterogeneous with special-purpose cores (e.g., multimedia, recognition, networking) and not only

GPUs mixed with general-purpose ones. Third, these different levels share many common problems such as efficient programming, scalability, and latency management. Hence, it is very important to conduct research targeting the heterogeneity at all presented hardware levels. Moreover, it is also important to take special care of the scalability issues, which form a key dimension in the complexity of today environment. The extreme scale of this environment comes from every level:

1. *Low Level*: number of CPUs, number of cores per processor;
2. *Medium Level*: number of nodes (e.g., with memory);
3. *High Level*: distributed/large-scale (geography dispersion, latency, etc.);
4. *Application*: extreme-scale problem size (e.g., calculation-intensive or data-intensive).

In 2008, the knowledge on how to efficiently use program or scale applications on such infrastructures was still vague. This was one of the main challenges that researchers wanted to take on. Therefore, at that time, we decided to launch the COST Action for high-performance and extreme-scale computing in such complex environments entitled “*Open European Network for High-Performance Computing in Complex Environments*.” The main reasons were as follows:

- There was a huge demand in terms of computational power for scientific and data-intensive applications;
- The architectural advances offered the potential to meet the application requirements;
- None of the state-of-the-art solutions in HPC at that time allowed exploitation to this potential level;
- Most of the research carried out in this area was fragmented and scattered across different research teams without any coordination.

COST¹ was indeed an appropriate framework for the proposed Action. The main goal of this Action was to overcome the actual research fragmentation on this very hot topic by gathering the most relevant European research teams involved in all the scientific areas described above (from the CPU core to the scientific applications) and coordinate their research.

Summarizing, this project within the COST framework allowed us to expect some potential benefits such as high-level scientific results in the very important domain of high-performance and extreme-scale computing in complex environment; strong coordination between different research teams with significant expertise on this subject; a better visibility of the European research in this area; and a strong impact on other scientists and high-performance applications.

¹European Cooperation in Science and Technology: <http://www.cost.eu>.

1.2 SCIENTIFIC ORGANIZATION

1.2.1 Scientific Focus

The expected scientific impacts of the project were to encourage the specific community to focus research on hot topics and applications of interest for the EU, to propagate the collaboration of research groups with the industry, to stimulate the formation of new groups in new EU countries, and to facilitate the solution of highly computationally demanding scientific problems as mentioned above. For this, the groups involved in this Action collaborated with several scientific and industrial groups that could benefit from the advances made by this Action, and prompted the incorporation of new groups to the network.

To achieve the research tasks, different leading European research teams participated in the concrete activities detailed in Section 1.3.

1.2.2 Working Groups

Four working groups were set up to coordinate the scientific research:

- numerical analysis for hierarchical and heterogeneous and multicore systems;
- libraries for the efficient use of complex systems with emphasis on computational library and communication library;
- algorithms and tools for mapping and executing applications onto distributed and heterogeneous systems;
- applications of hierarchical-heterogeneous systems.

It is important to note that these working groups targeted vertical aspects of the architectural structure outlined in the previous section. For instance, the Action's goal was to carry out work on numerical analysis at the multicore level, at the heterogeneous system level, as well as at the large-scale level. The last working group (Applications) was expected to benefit from research of the other three groups.

1.3 ACTIVITIES OF THE PROJECT

To achieve the goal of this Action, the following concrete activities were proposed. The main goal was to promote collaboration through science meetings, workshops, schools, and internships. This allowed interchange of ideas and mobility of researchers.

1.3.1 Spring Schools

The goal was to provide young researchers with a good opportunity to share information and knowledge and to present their current research. These schools contributed to the expansion of the computing community and spread of EU knowledge.

1.3.2 International Workshops

The goal of these meetings was to take the opportunity during international conferences to meet the attendees and other researchers by co-locating workshops.

1.3.3 Working Groups Meetings

The scientific work plan was divided among different working groups. Each working group had substantial autonomy in terms of research projects. A leader nominated by the Management Committee led each working group. Members of a given working group met once or twice a year to discuss and exchange specific scientific issues and problems.

1.3.4 Management Committee Meetings

These meetings were devoted to the organization of the network and ensured the scientific quality of the network.

1.3.5 Short-Term Scientific Missions

The goal of short-term scientific missions (STSMs) was to enable visits by early stage researchers to foreign laboratories and departments. This was mainly targeted at young researchers to receive cross-disciplinary training and to take advantage of the existing resources. The goal was to increase the competitiveness and career development of those scientists in this rapidly developing field through cutting-edge collaborative research on the topic.

1.4 MAIN OUTCOMES OF THE ACTION

We believe that this COST Action was a great success. It gathered 26 European countries and 2 non-COST countries (Russia and South Africa). We have held 12 meetings and 2 spring schools. Fifty-two STSMs have been carried out. We have a new FP7 project coming from this Action (HOST). We have edited a book, and more than 100 papers have been published thanks to this Action.

We have set up an application catalog that gathers applications from the Action members. Its goal is to gather a set of HPC applications that can be used as test cases or benchmarks for researchers in the HPC field. The applications catalog is available at <https://complexhpc-catalogue.bordeaux.inria.fr>.

In total, the Action gathered more than 250 participants over the four years of the project.

We have sent a survey to the Action members. From this survey, it clearly appears that one of the greatest successes of the Action is the continuous strengthening of the network for many of its members both in terms of research teams and research domains. Many STSMs have been done through new network connections. Spring schools are seen as a major success, as they helped many young researchers to share

and exchange knowledge and gain new connections. Many PhD theses have been defended during the course of the Action, and some of the management committee members have been invited on the defense board of some of these PhDs. Moreover, many presentations given during the meeting are considered very useful and have opened new research directions for other attendees.

We had four goals in this Action:

1. to train new generations of scientists in high-performance and heterogeneous computing;
2. to overcome research fragmentation, and foster HPC efforts to increase Europe's competitiveness;
3. to tackle the problem at every level (from cores to large-scale environment);
4. vertical integration to provide new integrated solutions for large-scale computing for future platforms.

Goal 1 has exceeded our expectations. The spring schools have been a great success. We had many STSMs, and the number of early stage researchers attending the meeting was always very high. We had great response from young researchers.

Goal 2 has also been achieved satisfactorily. Thanks to the Action, many joint researches have been carried out, and we have created a nice network of researchers within our Action. Moreover, many top-level publications have been made thanks to the Action.

Goal 3 has also been achieved. We have scientific results that cover the core level and the distributed infrastructure, as well as results that cover the intermediate layers. This is due to the fact that the consortium was made of researchers from different areas. This was very fruitful.

Goal 4 has not been achieved. The main reason is the fact that providing integrated solutions requires more research and development than a COST Action can provide. It goes far beyond the networking activities of COST Action.

1.5 CONTENTS OF THE BOOK

This book presents some of the main results, in terms of research, of the COST Action presented in this chapter. We are very proud to share this with the interested reader. We have structured the book according to the following parts in order to have a good balance between each part:

1. Numerical Analysis for Heterogeneous and Multicore Systems (Chapters 2, 3, and 4);
2. Communication and Storage Considerations in High-Performance Computing (Chapters 5, 6, 7, and 8);
3. Efficient Exploitation of Heterogeneous Architectures (Chapters 9, 10, 11, and 12);
4. CPU + GPU coprocessing (Chapters 13, 14, and 15);

5. Efficient Exploitation of Distributed Systems (Chapters 16 and 17);
6. Energy Awareness in High-Performance Computing (Chapters 18, 19, and 20);
7. Applications of Heterogeneous High-Performance Computing (Chapters 21, 22, and 23).

Chapter 2 discusses the redesign of the iterative solution algorithm in order to efficiently execute them on heterogeneous architectures. Chapter 3 studies the performance of a meshless numerical partial differential equation (PDE) solver, parallelized with OpenMP. The results depend on the way the computations are distributed and the way the cache is used. Chapter 4 presents the development of three parallel numerical algorithms for the solution of parabolic problems on graphs with a theoretical and experimental study of their scalability.

Chapter 5 surveys different techniques for mapping processes to computing units in order to optimize communication cost and reduce execution time. Chapter 6 offers a comprehensive overview of how to implement topology- and performance-aware collective communications. Chapter 7 analyzes the many-core architecture using a new model (K-model) in order to estimate the complexity of a given algorithm designed for such an architecture. Chapter 8 presents a scalable I/O storage system for the hierarchical architecture of Blue Gene computers featuring buffering and asynchronous I/O.

Chapter 9 describes algorithmic techniques for offline scheduling of independent workflows in order to satisfy user's quality of service. Chapter 10 investigates the advantage of using modern heterogeneous architecture for the efficient implementation of the Reed–Solomon erasure code. Chapter 11 analyzes the factors that enable the development of efficient parallel programs on modern many-core parallel architecture. Chapter 12 studies efficient solutions for electromagnetism applications in clusters of CPU + GPU nodes.

Chapter 13 describes how the functional performance model can be used to optimize the performance of scientific applications for heterogeneous and hierarchical platform. Chapter 14 presents algorithms for multilevel load-balancing on multicore and multi-GPU environments. Chapter 15 faces the all-pair shortest path problem for sparse graph. Different scheduling strategies are studied to efficiently solve such problems on heterogeneous systems.

Chapter 16 surveys different resource management systems and scheduling algorithms for HPC for clouds. Chapter 17 discusses different approaches for performing resource discovery in large-scale distributed systems.

Chapter 18 focuses on how to optimize and adapt software solution to improve energy efficiency in the context of HPC application. Chapter 19 studies energy-aware scheduling policies for three scenarios of federated cloud dealing with energy awareness. Chapter 20 explores the use of heterogeneous chip multiprocessors for network security and strategy to improve energy consumption in such contexts.

Chapter 21 describes the “jungle computing paradigm,” which consists in gathering a complex hierarchical collection of heterogeneous computing hardware with an application to hyperspectral remote sensing. Chapter 22 presents a new model for image and video processing based on parallel and heterogeneous platforms in order

to improve the performance of the application when dealing with high-definition images. Chapter 23 applies load-balancing techniques to efficiently execute tomographic reconstruction using hybrid GPU + CPU systems.

As you can see, this covers a large spectrum of results and topics on HPC and heterogeneous systems.

We wish you a fruitful and enjoyable time with this book.

ACKNOWLEDGMENT

This publication is supported by COST.

PART II

*Numerical Analysis for
Heterogeneous and
Multicore Systems*

2

On the Impact of the Heterogeneous Multicore and Many-Core Platforms on Iterative Solution Methods and Preconditioning Techniques

Dimitar Lukarski and Maya Neytcheva

Uppsala University, Uppsala, Sweden

Computer simulations are now broadly recognized as a third branch of research, complementing theory and experimental work. The significant increase of available computing power has enabled tackling very large scale, challenging, real-life problems and opening new possibilities for revolutionary breakthrough results in science and engineering. At the same time, the complexity of the computer architecture has risen to levels where it is possible to achieve its full computing power only after careful redesigning of existing algorithms and developing novel computational and

High-Performance Computing on Complex Environments, First Edition.

Edited by Emmanuel Jeannot and Julius Žilinskas.

© 2014 John Wiley & Sons, Inc. Published 2014 by John Wiley & Sons, Inc.

communication strategies. In this chapter, we discuss this issue for a class of methods, broadly used in scientific computations—the iterative solution methods.

2.1 INTRODUCTION

For many years, the potential of available, serial as well as parallel, computer resources has been growing hand in hand with the need to numerically solve increasingly larger models of real-life problems. During the past decades, it has been recognized that, together with theoretical development and laboratory or field experiments, computation has become a third branch of research. Scientific computing is today's driving force behind the progress in the most challenging and demanding problems we attempt to solve. As examples, we mention turbulent combustion with complex chemistry, atmospheric dynamics, laser fusion, medical imaging, detailed modeling of the human heart, and artificial brain simulation with over a million neurons, to name a few.

In recent years, we have been witnessing a change in the means to increase the computational power which has a strong impact on how that power should be utilized via the algorithms used in scientific computations. Therefore, we briefly describe the phases in performing numerical simulations. Consider a complex physical phenomenon, described as a set of, usually coupled, processes that develop in space and time, which we want to study, analyze, and predict. It is assumed that the simulation requires a large amount of computer resources in terms of memory and computation.

The process of performing the numerical simulations can be split into the following steps:

- I Mathematical model: Describes the phenomenon continuously in time and space in terms of mathematical relations, most often as coupled ordinary or partial differential equations. These equations depend on various problem parameters, such as thermal conductivity, capacitance, material properties, and so on.
- II Discrete model: Because of the high complexity of the continuous model, analytical solutions are in general not available. Therefore, we pose the task to compute the solution in a number of discrete points in time and space, thus discretizing the mathematical model. This can be accomplished using various techniques. Space discretization can be done using finite differences, finite elements, finite volumes, boundary elements, and so on. Similarly, in time, various explicit or implicit time-stepping procedures can be utilized. In addition to the model parameters, here additional discretization parameters are introduced, usually denoted as h in space and τ in time. The discrete model is expressed in terms of linear or nonlinear algebraic systems of equations which have to be solved. Depending on the problem, but also on the discretization techniques, the matrices associated with the algebraic systems can be dense or sparse, symmetric or nonsymmetric, and so on. As nonlinear systems are most

often solved via linearization, we consider from now on only linear systems that are also large and sparse.

- III The linear systems arising in Step II have to be solved by a proper solution method—direct or iterative. Because of the targeted large-sized problems and the lesser demands on computer resources, we consider iterative solvers only. The iterative methods may introduce yet other method parameters which increase further the dimension of the parameter space.
- IV Computer implementation: To enable computer simulations, the numerical methods have to be implemented on some computer platform.
- V Visualization and verification: This step is also of importance, but it is not considered here any further.

When performing numerical simulations, we deal with two major concerns. The first concern is *robustness* with respect to model, discretization, and method parameters. Robustness is understood in the sense that the numerical efficiency of the iterative method chosen in Step III should not depend on changes in the parameters. For example, the number of iterations should not increase uncontrollably when h decreases. The *numerical efficiency*, related to fast convergence rate, can be seen also as an element of the robustness of the method. The second concern is the *efficiency of the implementation* in Step IV. It is based on a programming model (such as shared or distributed memory model), programming language, and a particular computer platform.

It has been recognized that, in order to achieve fast, accurate, and reliable results from computer simulations, sufficient knowledge is required for all the above steps, in particular Steps II, III, and IV, and awareness of the interplay among them. By choosing one or another discretization method, we may influence the structure and the properties of the arising matrices; by choosing a particular solution method we may ensure robustness with respect to the various problem, discretization, and method parameters; but we may sacrifice the amount of internal parallelism. Knowledge about the computer architecture on which the simulations are to be run may influence the choice of the solution method, which in turn has to be combined with the requirements of accuracy and robustness.

With the ongoing radical shift in the computer architecture toward multicore and many-core computational units, the importance of the above arguments becomes even stronger. The new technology based on multicore and many-core devices provides higher performance capabilities both in terms of computational power (GFlop/s) and memory bandwidth (GB/s). The available and easily accessible power enables scientists to tackle larger problems with higher resolution, providing in this way a better understanding of the world.

The radical shift is clearly seen when we compare the supercomputers available 10 years ago with the personal computers of today. All supercomputers in 2003 contained less than 10,000 cores¹. By the end of 2004, the situation was still the

¹Top500 <http://www.top500.org/statistics/efficiency-power-cores/>.

same, with two exceptions. At present, the computer landscape is very different. Not only the Top500 leaders have over 500,000 cores. Currently, NVIDIA delivers GPU (graphical processing unit) cards with more than 2500 cores per device (see GPU NVIDIA K20X²). With an improved power supply, four of these cards can be installed in a standard personal computer and thus one can obtain a system with more than 10,000 cores, which is a commodity at our desktop.

In order to achieve fast and reliable performance of the iterative methods, it becomes crucial to reconsider and redesign the implementation of well-known algorithms as well as to gain a deeper insight into what the expected performance is of the most common solution techniques on multicore heterogeneous computer platforms.

The focus of this chapter is to show how iterative methods can be performed efficiently on highly parallel, heterogeneous platforms. We present various methods and examples to show how this can be done, which includes mathematical description, as well as hardware-specific aspects.

2.2 GENERAL DESCRIPTION OF ITERATIVE METHODS AND PRECONDITIONING

We briefly discuss basic iterative techniques as well as two of the most often used projection-based methods—the conjugate gradient (CG) method [1], the generalized minimal residual (GMRES) method [2], and the multigrid (MG) method [3]. We also describe the defect-correction technique [4] as an illustration of an approach particularly suitable for solving linear systems on heterogeneous computers.

2.2.1 Basic Iterative Methods

Consider the solution of the linear system

$$A\mathbf{x} = \mathbf{b}, \quad (2.1)$$

where $A \in \mathbb{R}^{n \times n}$ is a nonsingular matrix, so Equation (2.1) has a unique solution. The matrix A is large and sparse, and therefore the number of nonzero elements, $\text{nnz}(A)$, is proportional to the size of the matrix, n ; that is, $\text{nnz}(A) = O(n)$.

Finding the solution to Equation (2.1) is equivalent to finding the root of the equation

$$\mathbf{b} - A\mathbf{x} = \mathbf{0}. \quad (2.2)$$

One straightforward way to introduce simple iterative methods is to rewrite (2.2) as a fixed-point iteration, namely

$$\begin{aligned} &\text{for some given } \mathbf{x}^{(0)}, \text{ iterate} \\ &\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (\mathbf{b} - A\mathbf{x}^{(k)}), \quad k = 0, 1, \dots \text{ until convergence.} \end{aligned} \quad (2.3)$$

²NVIDIA K20 specification <http://www.nvidia.com/object/tesla-servers.html>

The computational procedure (2.3) defines a basic stationary iterative method. The computation cost per iteration involves one matrix–vector multiplication and two vector updates, and is clearly $O(n)$. Such a method, however, usually exhibits too slow a convergence, which manifests itself in unacceptably many iterations. In some cases, convergence may not even be achieved.

Aiming at accelerating the convergence of the iterative process has led to the idea to involve some method parameter, replacing the simple iteration (2.3) by

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau \mathbf{r}^{(k)}, \quad \text{or} \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau_k \mathbf{r}^{(k)}, \quad (2.4)$$

where $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$ is the *residual* at the k th iteration and τ or τ_k are some properly chosen method parameters. In Equation (2.4), the method parameters to tune are scalars. Of course, nothing prevents us from replacing them with a properly chosen matrix, referred to in the sequel as P ; thus we consider

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + P^{-1} \mathbf{r}^{(k)}, \quad k = 0, 1, \dots \text{ until convergence.} \quad (2.5)$$

As will be discussed later, P can also vary during the iterative process. For simplicity, now we consider that it is some explicitly given nonsingular matrix.

It is easy to see that Equation (2.5) is obtained by replacing the original system $A\mathbf{x} = \mathbf{b}$ by the transformed system

$$P^{-1}A\mathbf{x} = P^{-1}\mathbf{b},$$

and applying the fixed-point scheme to it. In this case, the iterative scheme becomes

$$\begin{cases} \mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}, \\ P\mathbf{d}^{(k)} = \mathbf{r}^{(k)}, \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}. \end{cases} \quad (2.6)$$

The scheme (2.6) has a higher computational complexity than that of Equation (2.4), since a solution of a system with the matrix P is required at each iteration. Clearly, the achieved decrease in the number of iterations must be significant enough to compensate for the extra cost per iteration.

We see from Equation (2.6) that the choice $P = A$ would lead to a procedure that converges within one iteration. However, the computational cost would be unacceptably high, similar to that of a direct solution method. Clearly, P should satisfy some conditions so that we can achieve faster convergence, keeping at the same time the overall computational costs of the whole iterative method as low as possible.

The matrix P is referred to as a *preconditioner* to A . We consider next some well-known choices of P , leading to a family of classical iterative methods which are based on the so-called matrix splitting technique.

Intuitively, P has to be related to A . Consider the following splitting of A ,

$$A = P - R,$$

where P is nonsingular and R can be seen as an error matrix. Then,

$$P^{-1}A = P^{-1}(P - R) = I - P^{-1}R,$$

where I is the identity matrix of proper order.

The matrix $B = P^{-1}R$ is referred to as the *iteration matrix* and is used in theoretical derivations to show the convergence of the corresponding iterative method, as well as to estimate its rate of convergence (see [5] for details). Using the splitting, we rewrite Equation (2.5) as follows:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + P^{-1}(\mathbf{b} - A\mathbf{x}^{(k)}) = P^{-1}\mathbf{b} + P^{-1}R\mathbf{x}^{(k)}$$

or

$$P\mathbf{x}^{(k+1)} = R\mathbf{x}^{(k)} + \mathbf{b}.$$

Let A be represented in the following way, $A = D - L - U$ where D , L , and U are the diagonal, the strictly lower triangular, and the strictly upper triangular part of A , respectively. Table 2.1 shows some classical iterative schemes, based on the latter splitting of A .

For more details on the convergence of these methods, refer to [5].

The common characteristic of these methods is the simplicity of their implementation. Here, P is a diagonal or a triangular matrix, and the degree of parallelism is related to the sparsity structure of the underlying matrices.

The bottleneck of these methods is their slow convergence. Their importance has not been lost, however. Today, they are mostly used as subsolvers in more advanced iterative techniques described in Sections 2.2.2 and 2.5. Because of their low arithmetic cost and ease of implementation, they are important ingredients in the so-called projection methods.

2.2.2 Projection Methods: CG and GMRES

The idea behind the projection methods is that the original problem of huge dimension (easily of tens or even hundreds of millions of degrees of freedom) is projected over a subspace of much smaller dimension. An approximate solution is sought in that smaller subspace and then projected back to the original large space. When the

TABLE 2.1 Classical Iterative Schemes Based on Matrix Splitting

Method	P	R	Scheme
Jacobi iteration	D	$L + U$	$D\mathbf{x}^{(k+1)} = (L + U)\mathbf{x}^{(k)} + \mathbf{b}$
Gauss–Seidel backward	$D - U$	L	$(D - U)\mathbf{x}^{(k+1)} = L\mathbf{x}^{(k)} + \mathbf{b}$
Gauss–Seidel forward	$D - L$	U	$(D - L)\mathbf{x}^{(k+1)} = U\mathbf{x}^{(k)} + \mathbf{b}$
SOR	$D - \omega L$	$\omega U + (1 - \omega)D$	$(D - \omega L)\mathbf{x}^{(k+1)} =$ $(\omega U + (1 - \omega)D)\mathbf{x}^{(k)} + \mathbf{b}$