
A SYSTOLIC ARRAY PARALLELIZING COMPILER

THE KLUWER INTERNATIONAL SERIES IN
ENGINEERING AND COMPUTER SCIENCE

HIGH PERFORMANCE COMPUTING
Systems, Networks and Algorithms

Consulting Editor

H.T. Kung

A SYSTOLIC ARRAY PARALLELIZING COMPILER

by

Ping-Sheng Tseng
Bell Communications Research, Inc.

with a foreword by

H.T. Kung



KLUWER ACADEMIC PUBLISHERS
Boston/Dordrecht/London

Distributors for North America:

Kluwer Academic Publishers
101 Philip Drive
Assinippi Park
Norwell, Massachusetts 02061 USA

Distributors for all other countries:

Kluwer Academic Publishers Group
Distribution Centre
Post Office Box 322
3300 AH Dordrecht, THE NETHERLANDS

Library of Congress Cataloging-in-Publication Data

Tseng, Ping-Sheng, 1959-

A systolic array parallelizing compiler / by Ping-Sheng Tseng.
p. cm. — (The Kluwer international series in engineering and
computer science ; #106. High performance computing)

Revision of the author's thesis (Ph.D.)—Carnegie Mellon
University.

Includes bibliographical references and index.

ISBN-13:978-1-4612-8835-0

e-ISBN-13:978-1-4613-1559-9

DOI: 10.1007/978-1-4613-1559-9

1. Parallel processing (Electronic computers) 2. Compilers
(Computer programs) I. Title. II. Series: Kluwer international
series in engineering and computer science ; SECS 106. III. Series:
Kluwer international series in engineering and computer science.
High performance computing.

QA76.58.T74 1990

005.4'53—dc20

90-4751

CIP

Copyright © 1990 by Kluwer Academic Publishers

Softcover reprint of the hardcover 1st edition 1990

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher, Kluwer Academic Publishers, 101 Philip Drive, Assinippi Park, Norwell, Massachusetts 02061.

Contents

Foreword	xi
Acknowledgements	xv
1 Introduction	1
2 Systolic array programming	5
2.1 The Warp machine	6
2.2 The W2 programming language	8
2.3 The AL programming language	11
2.4 Related work	18
3 Data relations	21
3.1 Linear data relations	24
3.2 Joint data compatibility classes	28
3.3 Scope of data compatibility classes	32
3.4 Summary	35
4 Loop Distribution	37
4.1 Intercell communication	38
4.2 The basic loop distribution scheme	41
4.3 Distributed loop parallelism	44
4.3.1 Intraloop parallelism	44
4.3.2 Interloop parallelism	45
4.4 Optimization	48
4.4.1 Load balancing	48

4.4.2	Communication scheduling	49
4.5	Related work	56
5	Implementation	59
5.1	External interface	59
5.2	Compiling DO* loops	61
5.3	The ALIGN* statement	62
5.4	Parallel accumulation	64
5.5	Program debugging	67
6	Evaluation	69
6.1	Matrix computations	73
6.1.1	LU decomposition	73
6.1.2	QR decomposition	81
6.1.3	Singular value decomposition	85
6.2	2D Fast Fourier Transform	89
6.3	Partial differential equation solvers	92
6.3.1	SOR	92
6.3.2	Line SOR	93
6.3.3	Two-color SOR	95
6.4	Summary	99
7	Conclusions	101
	Bibliography	105
A	Linear data relations in Livermore Loops	109
B	Benchmark programs	115
B.1	LU decomposition	115
B.1.1	Single cell	115
B.1.2	Multiple cells	116
B.2	QR decomposition	118
B.2.1	Single cell	118
B.2.2	Multiple cells	118
B.3	Singular value decomposition	120
B.3.1	Single cell	120

B.3.2	Multiple cell	121
B.4	2D Fast Fourier Transform	123
B.4.1	Single cell	123
B.4.2	Multiple cell	124
B.5	Partial differential equation solvers	125
B.5.1	SOR	125
B.5.2	Line SOR	126
B.5.3	Two-color SOR	127
Index		129

List of Figures

2.1	The Warp systolic array	6
2.2	The Warp cell architecture	7
2.3	A W2 matrix multiplication program	12
2.4	An AL matrix multiplication program	13
2.5	An AL compiler generated W2 program: matrix fac- torization	17
3.1	Maximal compatibility classes	22
4.1	Blocking	42
4.2	Interloop parallelism	47
4.3	Communication scheduling	52
4.4	An AL compiler generated W2 program: a relaxation loop	54
4.5	Communication scheduling: a relaxation loop	55
6.1	LU decomposition: 100×100 matrix	74
6.2	LU decomposition: 180×180 matrix	75
6.3	LU decomposition with parallel pivoting: 100×100 matrix	79
6.4	LU decomposition with parallel pivoting: 180×180 matrix	80
6.5	QR decomposition: 100×100 matrix	83
6.6	QR decomposition: 180×180 matrix	84
6.7	SVD bidiagonalization: 100×100 matrix	87
6.8	SVD bidiagonalization: 175×175 matrix	88
6.9	2D FFT, 64×64 complex image	90

6.10 SOR 94

6.11 Line SOR 96

6.12 Two-color SOR 98

Foreword

Widespread use of parallel processing will become a reality only if the process of porting applications to parallel computers can be largely automated. Usually it is straightforward for a user to determine how an application can be mapped onto a parallel machine; however, the actual development of parallel code, if done by hand, is typically difficult and time consuming. Parallelizing compilers, which can generate parallel code automatically, are therefore a key technology for parallel processing.

In this book, Ping-Sheng Tseng describes a parallelizing compiler for systolic arrays, called AL. Although parallelizing compilers are quite common for shared-memory parallel machines, the AL compiler is one of the first working parallelizing compilers for distributed-memory machines, of which systolic arrays are a special case. The AL compiler takes advantage of the fine grain and high bandwidth interprocessor communication capabilities in a systolic architecture to generate efficient parallel code.

While capable of handling an important class of applications, AL is not intended to be a general-purpose parallelizing compiler. Instead, AL is designed to be effective for a special class of computations that use arrays and loops. AL relies on the fact that for these computations, the user can easily provide “hints” or mapping strategies to guide the compiler to distribute data structures and loop iterations. Using these hints, the AL compiler generates the local program for each processor, manages the interprocessor communication, and parallelizes the loop execution. A fundamental contribution of AL, which goes beyond its current implementation, is the identification of what to capture in these hints that the user can easily provide and the compiler can effectively use.

AL has proven to be extremely effective in programming the Warp systolic array developed by Carnegie Mellon. AL was used to port the nontrivial LINPACK QR (SQRDC), SVD (SSVDC), LU (SGEFA), and back substitution (SGESL) routines to Warp all in one person-week. AL has been used in several applications, including the porting of a large Navy signal processing application to Warp. The AL compiler has also been used to generate parallel code for automatic schedulers that map a large set of high-level signal processing tasks onto Warp.

Carnegie Mellon’s experience in programming Warp clearly indi-

cates the effectiveness of special-purpose parallelizing compilers such as AL. These tools, which we also call parallel program generators, can improve a programmer's productivity by several orders of magnitude. In fact, some applications would never have been brought up on Warp if the user did not have access to AL; for these applications explicit programming of interprocessor communication is simply too difficult to be done by a human being. Besides AL, the Warp project has developed several other parallel program generators. One of them, called Apply, has been extensively used to generate parallel code for image processing applications. Both AL and Apply generally produce code better than or as good as hand-written code.

The success of parallel program generators such as AL is encouraging. With them we are assured that programming parallel machines can be as easy as programming sequential machines for some important application areas. In this sense, these programming tools have helped legitimize our effort in building parallel computers. Besides providing useful tools for programming parallel machines for special applications, these parallel program generators are also significant in giving insights into information that more general-purpose parallelizing compilers of the future need to capture.

The book is an outgrowth of Ping-Sheng Tseng's Ph.D thesis from Carnegie Mellon. The book gives a complete treatment of the

design, implementation and evaluation of AL. I am pleased to write the Foreword for this outstanding piece of work and hope that the book will inspire researchers to further this very important area of parallel processing.

H. T. Kung
April, 1990
Pittsburgh, PA

Acknowledgements

This work evolved from my thesis research at Carnegie Mellon University. I want to thank my thesis advisor H. T. Kung. When I started this research project, Kung was the only one who believed that I could make this thesis happen. Without his vision, support and encouragement, I might have changed my thesis topic before I started. Thanks to Monica Lam and Peter Lieu. Monica implemented the W2 compiler directives for me. Peter has been very responsive to my W2 bug reports. Without their help, all the experimental work might have never been completed. The members of my thesis committee are H. T. Kung, Gary Koob, Thomas Gross, and George Cox. Their invaluable comments helped me to improve the quality of this presentation. I want to thank my former officemate Robert Cohn. Robert read my manuscript several times and helped me to make this presentation much more readable. I would also like to thank the Defense Advanced Research Projects Agency for supporting this work. Finally, I would like to thank my wife Ly-June

and my father Po and my mother Sou-Sou for their support in my graduate study.

A SYSTOLIC ARRAY PARALLELIZING COMPILER