

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

## MICROPROCESSOR ARCHITECTURE

This book gives a comprehensive description of the architecture of microprocessors from simple in-order short pipeline designs to out-of-order superscalars. It discusses topics such as

- The policies and mechanisms needed for out-of-order processing, such as register renaming, reservation stations, and reorder buffers
- Optimizations for high performance, such as branch predictors, instruction scheduling, and load-store speculations
- Design choices and enhancements to tolerate latency in the cache hierarchy of single and multiple processors
- State-of-the-art multithreading and multiprocessing, emphasizing single-chip implementations

Topics are presented as conceptual ideas, with metrics to assess the effects on performance, if appropriate, and examples of realization. The emphasis is on how things work at a black box and algorithmic level. The author also provides sufficient detail at the register transfer level so that readers can appreciate how design features enhance performance as well as complexity.

Jean-Loup Baer is Professor Emeritus of Computer Science and Engineering at the University of Washington, where he has been since 1969. Professor Baer is the author of *Computer Systems Architecture* and of more than 100 refereed papers. He is a Guggenheim Fellow, an ACM Fellow, and an IEEE Fellow. Baer has held several editorial positions, including editor-in-chief of the *Journal of VLSI and Computer Systems* and editor of the *IEEE Transactions on Computers*, the *IEEE Transactions on Parallel and Distributed Systems*, and the *Journal of Parallel and Distributed Computing*. He has served as General Chair and Program Chair of several conferences, including ISCA and HPCA.

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

---

# Microprocessor Architecture

## FROM SIMPLE PIPELINES TO CHIP MULTIPROCESSORS

**Jean-Loup Baer**

University of Washington, Seattle



**CAMBRIDGE**  
UNIVERSITY PRESS

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

CAMBRIDGE UNIVERSITY PRESS

Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore,  
São Paulo, Delhi, Dubai, Tokyo

Cambridge University Press

32 Avenue of the Americas, New York, NY 10013-2473, USA

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9780521769921](http://www.cambridge.org/9780521769921)

© Jean-Loup Baer 2010

This publication is in copyright. Subject to statutory exception  
and to the provisions of relevant collective licensing agreements,  
no reproduction of any part may take place without the written  
permission of Cambridge University Press.

First published 2010

Printed in the United States of America

*A catalog record for this publication is available from the British Library.*

*Library of Congress Cataloging in Publication data*

Baer, Jean-Loup.

Microprocessor architecture: from simple pipelines to chip multiprocessors /  
Jean-Loup Baer.

p. cm.

Includes bibliographical references and index.

ISBN 978-0-521-76992-1 (hardback)

1. Microprocessors. 2. Computer architecture. I. Title.

QA76.5.B227 2009

004.2'2 – dc22 2009025686

ISBN 978-0-521-76992-1 Hardback

Cambridge University Press has no responsibility for the persistence or  
accuracy of URLs for external or third-party Internet Web sites referred to in  
this publication and does not guarantee that any content on such Web sites is,  
or will remain, accurate or appropriate.

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

---

*To Diane, Marc, Shawn, and Danielle*

## Contents

<i>Preface</i>	<i>page xi</i>
<b>1 Introduction</b> . . . . .	<b>1</b>
1.1 A Quick View of Technological Advances	2
1.2 Performance Metrics	6
1.3 Performance Evaluation	12
1.4 Summary	22
1.5 Further Reading and Bibliographical Notes	23
EXERCISES	24
REFERENCES	28
<b>2 The Basics</b> . . . . .	<b>29</b>
2.1 Pipelining	29
2.2 Caches	46
2.3 Virtual Memory and Paging	59
2.4 Summary	68
2.5 Further Reading and Bibliographical Notes	68
EXERCISES	69
REFERENCES	73
<b>3 Superscalar Processors</b> . . . . .	<b>75</b>
3.1 From Scalar to Superscalar Processors	75
3.2 Overview of the Instruction Pipeline of the DEC Alpha 21164	78
3.3 Introducing Register Renaming, Reorder Buffer, and Reservation Stations	89
3.4 Overview of the Pentium P6 Microarchitecture	102
3.5 VLIW/EPIC Processors	111
3.6 Summary	121
3.7 Further Reading and Bibliographical Notes	122
EXERCISES	124
REFERENCES	126

<b>4 Front-End: Branch Prediction, Instruction Fetching, and Register Renaming</b> . . . . .	129
4.1 Branch Prediction	130
Sidebar: The DEC Alpha 21264 Branch Predictor	157
4.2 Instruction Fetching	158
4.3 Decoding	164
4.4 Register Renaming (a Second Look)	165
4.5 Summary	170
4.6 Further Reading and Bibliographical Notes	170
EXERCISES	171
Programming Projects	174
REFERENCES	174
<b>5 Back-End: Instruction Scheduling, Memory Access Instructions, and Clusters</b> . . . . .	177
5.1 Instruction Issue and Scheduling (Wakeup and Select)	178
5.2 Memory-Accessing Instructions	184
5.3 Back-End Optimizations	195
5.4 Summary	203
5.5 Further Reading and Bibliographical Notes	204
EXERCISES	205
Programming Project	206
REFERENCES	206
<b>6 The Cache Hierarchy</b> . . . . .	208
6.1 Improving Access to L1 Caches	209
6.2 Hiding Memory Latencies	218
6.3 Design Issues for Large Higher-Level Caches	232
6.4 Main Memory	245
6.5 Summary	253
6.6 Further Reading and Bibliographical Notes	254
EXERCISES	255
Programming Projects	257
REFERENCES	258
<b>7 Multiprocessors</b> . . . . .	260
7.1 Multiprocessor Organization	261
7.2 Cache Coherence	269
7.3 Synchronization	281
7.4 Relaxed Memory Models	290
7.5 Multimedia Instruction Set Extensions	294
7.6 Summary	296
7.7 Further Reading and Bibliographical Notes	297
EXERCISES	298
REFERENCES	300

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

## Contents

ix

<b>8 Multithreading and (Chip) Multiprocessing</b> . . . . .	303
8.1 Single-Processor Multithreading	304
8.2 General-Purpose Multithreaded Chip Multiprocessors	318
8.3 Special-Purpose Multithreaded Chip Multiprocessors	324
8.4 Summary	330
8.5 Further Reading and Bibliographical Notes	331
EXERCISES	332
REFERENCES	333
<b>9 Current Limitations and Future Challenges</b> . . . . .	335
9.1 Power and Thermal Management	336
9.2 Technological Limitations: Wire Delays and Pipeline Depths	343
9.3 Challenges for Chip Multiprocessors	346
9.4 Summary	348
9.5 Further Reading and Bibliographical Notes	349
REFERENCES	349
<i>Bibliography</i>	351
<i>Index</i>	361

## Preface

Computer architecture is at a turning point. Radical changes occurred in the 1980s when the Reduced Instruction Set Computer (RISC) philosophy, spurred in good part by academic research, permeated the industry as a reaction to the Complex Instruction Set Computer (CISC) complexities. Today, three decades later, we have reached a point where physical limitations such as power dissipation and design complexity limit the speed and the performance of single-processor systems. The era of chip multiprocessors (CMP), or multicores, has arrived.

Besides the uncertainty on the structure of CMPs, it is not known yet whether the future CMPs will be composed of simple or complex processors or a combination of both and how the on-chip and off-chip memory hierarchy will be managed. It is important for computer scientists and engineers to look at the possible options for the modules that will compose the new generations of multicores. In this book, we describe the architecture of microprocessors from simple in-order short pipe designs to out-of-order superscalars with many optimizations. We also present choices and enhancements in the cache hierarchy of single processors. The last part of this book introduces readers to the state-of-the-art in multithreading and multiprocessing, emphasizing single-chip implementations and their cache hierarchy.

The emphasis in this book is on “how things work” at a black box and algorithmic level. It is not as close to the hardware as books that explain features at the register transfer level. However, it provides sufficient detail, say at the “bit level” and through pseudocode, so that the reader can appreciate design features that have or will enhance performance as well as their complexity.

As much as possible we present topics as conceptual ideas, define metrics to assess the performance impact if appropriate, show alternate (if any) ways of implementing them, and provide examples of realization.

### Synopsis

This book has three main parts. Chapter 1 and Chapter 2 review material that should have been taught in a prerequisite class. Chapters 3 through 6 describe single-processor systems and their memory hierarchy. Chapter 7 and Chapter 8 are



devoted to parallelism. The last (short) chapter, Chapter 9, introduces limitations due to technology and presents challenges for future CMPs.

More specifically:

- Chapter 1 reviews Moore's law and its influence as well as current limitations to ever faster processors. Performance metrics such as *cycles per instruction* (CPI), *instruction per cycle* (IPC), and speedup are defined. A section on performance evaluation introduces benchmarks and simulators. Chapter 2 summarizes the main concepts behind pipelining, including data and control hazards. Two versions of the "classical" 5-stage pipeline are contrasted. The basics of cache organization and performance are reviewed. Virtual memory (paging) and *Translation Look-Aside Buffers* (TLBs) are presented from the computer architecture (not the operating system) viewpoint.
- Chapter 3 describes two landmark "families" in the history of modern multiple-issue microprocessors: in-order processors represented by the DEC Alpha 21164 and out-of-order processors represented by the Intel Pentium P6 micro-architecture. In order to ease the description of the latter, a first look at register renaming is provided, with a sidebar showing the classical Tomasulo algorithm. This chapter ends with an introduction to the Very Long Instruction Word (VLIW) / Explicitly Parallel Instruction Computing (EPIC) philosophy and a brief overview of its most well-known industrial realization, the Intel Itanium. Chapter 4 and Chapter 5 are devoted to detailed explanations of, respectively, the front end and the back end of the pipeline. Because many advances in superscalar performance have been attained through speculation, we start Chapter 4 with a "model" of a branch predictor. This model will be referred to in the book when other speculative schemes are explored. We then describe actual branch predictors and branch target predictors. A sidebar describes the sophisticated Alpha 21264 tournament predictor. The remainder of the chapter deals with instruction fetching, including trace caches, and presents another view of register renaming that has become more popular than Tomasulo's. Chapter 5 looks at another potential bottleneck in the pipeline, namely the wakeup-select cycle. Load speculation and other back-end optimizations are introduced. Chapter 6 deals with the cache hierarchy: how to improve access times for first-level caches (another critical potential source of slowdown); methods to hide latency from higher-level caches and main memory such as prefetching and lock-up free caches; and, in a more research-oriented way, design and performance issues for large caches. We conclude this chapter with a look at main memory.
- Chapter 7 and Chapter 8 present multiprocessors and multithreading. In Chapter 7, we first introduce taxonomies that will set the stage for the two chapters. Cache coherence is then treated with both snoopy protocols and directory protocols, because it is likely that both types of protocols will be used in future parallel processing systems, whether on-chip or between chips. The next topic is synchronization, in which, in addition to the existing lock mechanisms, we mention approaches that might replace or complement them, such as transactional

memory. We also briefly introduce relaxed memory models, a topic that may become more important in the near future. The chapter concludes with a description of multimedia instruction set extensions because this is a (limited) form of parallelism. In Chapter 8, we start by looking at various flavors of multithreading: fine-grained, coarse-grained, and SMT. This leads us to the description of some current CMPs. We start with two examples of general-purpose CMPs: the Sun Niagara, a CMP using fine-grained multithreaded simple processors, and the Intel Core Duo, a representative of Intel multicores using complex superscalar processors. We conclude with two special-purpose CMPs: the IBM Cell, intended for games and scientific computation, and the Intel IXP, a network processor using coarse-grained multithreaded microengines.

- Finally, Chapter 9 gives a brief introduction to the factors that have capped the exponential performance curve of single processors, namely power issues, wire lengths, and pipeline depths. We end the chapter with an (incomplete) list of challenges for future CMPs.

### Use of the book

This book grew out of courses given over the last ten years at the senior undergraduate and first-year graduate level for computer science and computer engineering students at the University of Washington. It is intended as a book for a second course in computer architecture. In an undergraduate class, we typically spend 10% of the time on review and presentation of simulators. We spend 60% of the time on single processors. About 25% is spent on multiprocessors and multithreading, and the remainder on some research presentations by colleagues and advanced graduate students. At the graduate level, the choice of topics depends more on the instructor. For example, one possible approach would be to assign the first two chapters as required reading and start directly with topics of the instructor's choice. A rule of thumb that we have followed is to have a 50–50 split in classroom time between single-processor (Chapters 3 through 6) and CMPs (Chapters 7 through 9).

### Acknowledgments

I am indebted to students who took my classes and to many colleagues, including my doctoral students, for their help and comments. I want especially to recognize the leaders and participants of the weekly CSE “Computer Architecture Lunch,” which Wen-Hann Wang and I created more than two decades ago. Their choice of the topics in this seminar and their constructive criticism of the importance of specific areas have influenced the contents of this book. Their patient explanations of what were for me obscure points has hopefully led to clearer exposition. In particular, I thank my faculty colleagues Luis Ceze, Carl Ebeling, Susan Eggers, Mark Oskin, and Larry Snyder and my research collaborators Craig Anderson, Tien-Fu Chen, Patrick Crowley, Dennis Lee, Douglas Low, Sang-Lyul Min, Xiaohan Qin, Taylor VanVleet, Wen-Hann Wang, Wayne Wong, and Rick Zucker.

Cambridge University Press

978-0-521-76992-1 - Microprocessor Architecture: From Simple Pipelines to Chip Multiprocessors

Jean-Loup Baer

Frontmatter

[More information](#)

---

xiv

Preface

I thank Lauren Cowles and David Jou at Cambridge University Press for providing me with excellent anonymous reviews, and Shana Meyer and colleagues at Aptara, Inc. – in particular, copy editor Joseph C. Fineman – for a smooth and timely production process.

Jean-Loup Baer  
Seattle, WA  
August 2009