
HIGH PERFORMANCE SWITCHES AND ROUTERS

H. JONATHAN CHAO and BIN LIU



WILEY-INTERSCIENCE
A JOHN WILEY & SONS, INC., PUBLICATION

HIGH PERFORMANCE SWITCHES AND ROUTERS



THE WILEY BICENTENNIAL—KNOWLEDGE FOR GENERATIONS

Each generation has its unique needs and aspirations. When Charles Wiley first opened his small printing shop in lower Manhattan in 1807, it was a generation of boundless potential searching for an identity. And we were there, helping to define a new American literary tradition. Over half a century later, in the midst of the Second Industrial Revolution, it was a generation focused on building the future. Once again, we were there, supplying the critical scientific, technical, and engineering knowledge that helped frame the world. Throughout the 20th Century, and into the new millennium, nations began to reach out beyond their own borders and a new international community was born. Wiley was there, expanding its operations around the world to enable a global exchange of ideas, opinions, and know-how.

For 200 years, Wiley has been an integral part of each generation's journey, enabling the flow of information and understanding necessary to meet their needs and fulfill their aspirations. Today, bold new technologies are changing the way we live and learn. Wiley will be there, providing you the must-have knowledge you need to imagine new worlds, new possibilities, and new opportunities.

Generations come and go, but you can always count on Wiley to provide you the knowledge you need, when and where you need it!

WILLIAM J. PESCE
PRESIDENT AND CHIEF EXECUTIVE OFFICER

PETER BOOTH WILEY
CHAIRMAN OF THE BOARD

HIGH PERFORMANCE SWITCHES AND ROUTERS

H. JONATHAN CHAO and BIN LIU



WILEY-INTERSCIENCE
A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2007 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) 750-4470, 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, or online at <http://www.wiley.com/go/permission>.

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 or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States 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 may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data.

Chao, H. Jonathan, 1955-

High performance switches and routers / by H. Jonathan Chao, Bin Liu.
p. cm.

ISBN-13: 978-0-470-05367-6

ISBN-10: 0-470-05367-4

1. Asynchronous transfer mode. 2. Routers (Computer networks)
3. Computer network protocols. 4. Packet switching (Data transmission)

I. Liu, Bin. II. Title.

TK5105.35.C454 2007

621.382'16--dc22

2006026971

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

CONTENTS

PREFACE	xv
ACKNOWLEDGMENTS	xvii
1 INTRODUCTION	1
1.1 Architecture of the Internet: Present and Future / 2	
1.1.1 The Present / 2	
1.1.2 The Future / 4	
1.2 Router Architectures / 5	
1.3 Commercial Core Router Examples / 9	
1.3.1 T640 TX-Matrix / 9	
1.3.2 Carrier Routing System (CRS-1) / 11	
1.4 Design of Core Routers / 13	
1.5 IP Network Management / 16	
1.5.1 Network Management System Functionalities / 16	
1.5.2 NMS Architecture / 17	
1.5.3 Element Management System / 18	
1.6 Outline of the Book / 19	
2 IP ADDRESS LOOKUP	25
2.1 Overview / 25	
2.2 Trie-Based Algorithms / 29	
2.2.1 Binary Trie / 29	
2.2.2 Path-Compressed Trie / 31	

2.2.3	Multi-Bit Trie / 33
2.2.4	Level Compression Trie / 35
2.2.5	Lulea Algorithm / 37
2.2.6	Tree Bitmap Algorithm / 42
2.2.7	Tree-Based Pipelined Search / 45
2.2.8	Binary Search on Prefix Lengths / 47
2.2.9	Binary Search on Prefix Range / 48
2.3	Hardware-Based Schemes / 51
2.3.1	DIR-24-8-BASIC Scheme / 51
2.3.2	DIR-Based Scheme with Bitmap Compression (BC-16-16) / 53
2.3.3	Ternary CAM for Route Lookup / 57
2.3.4	Two Algorithms for Reducing TCAM Entries / 58
2.3.5	Reducing TCAM Power – CoolCAMs / 60
2.3.6	TCAM-Based Distributed Parallel Lookup / 64
2.4	IPv6 Lookup / 67
2.4.1	Characteristics of IPv6 Lookup / 67
2.4.2	A Folded Method for Saving TCAM Storage / 67
2.4.3	IPv6 Lookup via Variable-Stride Path and Bitmap Compression / 69
2.5	Comparison / 73

3 PACKET CLASSIFICATION

77

3.1	Introduction / 77
3.2	Trie-Based Classifications / 81
3.2.1	Hierarchical Tries / 81
3.2.2	Set-Pruning Trie / 82
3.2.3	Grid of Tries / 83
3.2.4	Extending Two-Dimensional Schemes / 84
3.2.5	Field-Level Trie Classification (FLTC) / 85
3.3	Geometric Algorithms / 90
3.3.1	Background / 90
3.3.2	Cross-Producing Scheme / 91
3.3.3	Bitmap-Intersection / 92
3.3.4	Parallel Packet Classification (P^2C) / 93
3.3.5	Area-Based Quadtree / 95
3.3.6	Hierarchical Intelligent Cuttings / 97
3.3.7	HyperCuts / 98
3.4	Heuristic Algorithms / 103
3.4.1	Recursive Flow Classification / 103
3.4.2	Tuple Space Search / 107

- 3.5 TCAM-Based Algorithms / 108
 - 3.5.1 Range Matching in TCAM-Based Packet Classification / 108
 - 3.5.2 Range Mapping in TCAMs / 110

4 TRAFFIC MANAGEMENT

114

- 4.1 Quality of Service / 114
 - 4.1.1 QoS Parameters / 115
 - 4.1.2 Traffic Parameters / 116
- 4.2 Integrated Services / 117
 - 4.2.1 Integrated Service Classes / 117
 - 4.2.2 IntServ Architecture / 117
 - 4.2.3 Resource ReSerVation Protocol (RSVP) / 119
- 4.3 Differentiated Services / 121
 - 4.3.1 Service Level Agreement / 122
 - 4.3.2 Traffic Conditioning Agreement / 123
 - 4.3.3 Differentiated Services Network Architecture / 123
 - 4.3.4 Network Boundary Traffic Classification and Conditioning / 124
 - 4.3.5 Per Hop Behavior (PHB) / 126
 - 4.3.6 Differentiated Services Field / 127
 - 4.3.7 PHB Implementation with Packet Schedulers / 128
- 4.4 Traffic Policing and Shaping / 129
 - 4.4.1 Location of Policing and Shaping Functions / 130
 - 4.4.2 ATM's Leaky Bucket / 131
 - 4.4.3 IP's Token Bucket / 133
 - 4.4.4 Traffic Policing / 134
 - 4.4.5 Traffic Shaping / 135
- 4.5 Packet Scheduling / 136
 - 4.5.1 Max-Min Scheduling / 136
 - 4.5.2 Round-Robin Service / 138
 - 4.5.3 Weighted Round-Robin Service / 139
 - 4.5.4 Deficit Round-Robin Service / 140
 - 4.5.5 Generalized Processor Sharing (GPS) / 141
 - 4.5.6 Weighted Fair Queuing (WFQ) / 146
 - 4.5.7 Virtual Clock / 150
 - 4.5.8 Self-Clocked Fair Queuing / 153
 - 4.5.9 Worst-Case Fair Weighted Fair Queuing (WF²Q) / 155
 - 4.5.10 WF²Q+ / 158
 - 4.5.11 Comparison / 159
 - 4.5.12 Priorities Sorting Using a Sequencer / 160

- 4.6 Buffer Management / 163
 - 4.6.1 Tail Drop / 163
 - 4.6.2 Drop on Full / 164
 - 4.6.3 Random Early Detection (RED) / 164
 - 4.6.4 Differential Dropping: RIO / 167
 - 4.6.5 Fair Random Early Detection (FRED) / 168
 - 4.6.6 Stabilized Random Early Detection (SRED) / 170
 - 4.6.7 Longest Queue Drop (LQD) / 172

5 BASICS OF PACKET SWITCHING 176

- 5.1 Fundamental Switching Concept / 177
- 5.2 Switch Fabric Classification / 181
 - 5.2.1 Time-Division Switching / 181
 - 5.2.2 Space-Division Switching / 183
- 5.3 Buffering Strategy in Switching Fabrics / 187
 - 5.3.1 Shared-Memory Queuing / 188
 - 5.3.2 Output Queuing (OQ) / 188
 - 5.3.3 Input Queuing / 189
 - 5.3.4 Virtual Output Queuing (VOQ) / 189
 - 5.3.5 Combined Input and Output Queuing / 190
 - 5.3.6 Crosspoint Queuing / 191
- 5.4 Multiplane Switching and Multistage Switching / 191
- 5.5 Performance of Basic Switches / 195
 - 5.5.1 Traffic Model / 196
 - 5.5.2 Input-Buffered Switches / 197
 - 5.5.3 Output-Buffered Switches / 199
 - 5.5.4 Completely Shared-Buffered Switches / 201

6 SHARED-MEMORY SWITCHES 207

- 6.1 Linked List Approach / 208
- 6.2 Content Addressable Memory Approach / 213
- 6.3 Space-Time-Space Approach / 215
- 6.4 Scaling the Shared-Memory Switches / 217
 - 6.4.1 Washington University Gigabit Switch / 217
 - 6.4.2 Concentrator-Based Growable Switch Architecture / 218
 - 6.4.3 Parallel Shared-Memory Switches / 218
- 6.5 Multicast Shared-Memory Switches / 220
 - 6.5.1 Shared-Memory Switch with a Multicast Logical Queue / 220
 - 6.5.2 Shared-Memory Switch with Cell Copy / 220
 - 6.5.3 Shared-Memory Switch with Address Copy / 222

7	INPUT-BUFFERED SWITCHES	225
7.1	Scheduling in VOQ-Based Switches /	226
7.2	Maximum Matching /	229
7.2.1	Maximum Weight Matching /	229
7.2.2	Approximate MWM /	229
7.2.3	Maximum Size Matching /	230
7.3	Maximal Matching /	231
7.3.1	Parallel Iterative Matching (PIM) /	232
7.3.2	Iterative Round-Robin Matching (<i>i</i> RRM) /	233
7.3.3	Iterative Round-Robin with SLIP (<i>i</i> SLIP) /	234
7.3.4	FIRM /	241
7.3.5	Dual Round-Robin Matching (DRRM) /	241
7.3.6	Pipelined Maximal Matching /	245
7.3.7	Exhaustive Dual Round-Robin Matching (EDRRM) /	248
7.4	Randomized Matching Algorithms /	249
7.4.1	Randomized Algorithm with Memory /	250
7.4.2	A Derandomized Algorithm with Memory /	250
7.4.3	Variant Randomize Matching Algorithms /	251
7.4.4	Polling Based Matching Algorithms /	254
7.4.5	Simulated Performance /	258
7.5	Frame-based Matching /	262
7.5.1	Reducing the Reconfiguration Frequency /	263
7.5.2	Fixed Size Synchronous Frame-Based Matching /	267
7.5.3	Asynchronous Variable-Size Frame-Based Matching /	270
7.6	Stable Matching with Speedup /	273
7.6.1	Output-Queuing Emulation with Speedup of 4 /	274
7.6.2	Output-Queuing Emulation with Speedup of 2 /	275
7.6.3	Lowest Output Occupancy Cell First (LOOFA) /	278
8	BANYAN-BASED SWITCHES	284
8.1	Banyan Networks /	284
8.2	Batcher-Sorting Network /	287
8.3	Output Contention Resolution Algorithms /	288
8.3.1	Three-Phase Implementation /	288
8.3.2	Ring Reservation /	288
8.4	The Sunshine Switch /	292
8.5	Deflection Routing /	294
8.5.1	Tandem Banyan Switch /	294
8.5.2	Shuffle-Exchange Network with Deflection Routing /	296
8.5.3	Dual Shuffle-Exchange Network with Error-Correcting Routing /	297

- 8.6 Multicast Copy Networks / 303
 - 8.6.1 Broadcast Banyan Network / 304
 - 8.6.2 Encoding Process / 308
 - 8.6.3 Concentration / 309
 - 8.6.4 Decoding Process / 310
 - 8.6.5 Overflow and Call Splitting / 310
 - 8.6.6 Overflow and Input Fairness / 311

9 KNOCKOUT-BASED SWITCHES 316

- 9.1 Single-Stage Knockout Switch / 317
 - 9.1.1 Basic Architecture / 317
 - 9.1.2 Knockout Concentration Principle / 318
 - 9.1.3 Construction of the Concentrator / 320
- 9.2 Channel Grouping Principle / 323
 - 9.2.1 Maximum Throughput / 324
 - 9.2.2 Generalized Knockout Principle / 325
- 9.3 Two-Stage Multicast Output-Buffered ATM Switch (MOBAS) / 327
 - 9.3.1 Two-Stage Configuration / 327
 - 9.3.2 Multicast Grouping Network (MGN) / 330
- 9.4 Appendix / 333

10 THE ABACUS SWITCH 336

- 10.1 Basic Architecture / 337
- 10.2 Multicast Contention Resolution Algorithm / 340
- 10.3 Implementation of Input Port Controller / 342
- 10.4 Performance / 344
 - 10.4.1 Maximum Throughput / 344
 - 10.4.2 Average Delay / 347
 - 10.4.3 Cell Loss Probability / 349
- 10.5 ATM Routing and Concentration (ARC) Chip / 351
- 10.6 Enhanced Abacus Switch / 354
 - 10.6.1 Memoryless Multi-Stage Concentration Network / 354
 - 10.6.2 Buffered Multi-Stage Concentration Network / 357
 - 10.6.3 Resequencing Cells / 359
 - 10.6.4 Complexity Comparison / 361
- 10.7 Abacus Switch for Packet Switching / 362
 - 10.7.1 Packet Interleaving / 362
 - 10.7.2 Cell Interleaving / 364

11 CROSSPOINT BUFFERED SWITCHES 367

- 11.1 Combined Input and Crosspoint Buffered Switches / 368

- 11.2 Combined Input and Crosspoint Buffered Switches with VOQ / 370
 - 11.2.1 CIXB with One-Cell Crosspoint Buffers (CIXB-1) / 371
 - 11.2.2 Throughput and Delay Performance / 371
 - 11.2.3 Non-Negligible Round-Trip Times in CIXB- k / 376
- 11.3 OCF_OCF: Oldest Cell First Scheduling / 376
- 11.4 LQF_RR: Longest Queue First and Round-Robin Scheduling in CIXB-1 / 378
- 11.5 MCBF: Most Critical Buffer First Scheduling / 379

12 CLOS-NETWORK SWITCHES

382

- 12.1 Routing Property of Clos Network Switches / 383
- 12.2 Looping Algorithm / 387
- 12.3 m -Matching Algorithm / 388
- 12.4 Euler Partition Algorithm / 388
- 12.5 Karol's Algorithm / 389
- 12.6 Frame-Based Matching Algorithm for Clos Network (f-MAC) / 391
- 12.7 Concurrent Matching Algorithm for Clos Network (c-MAC) / 392
- 12.8 Dual-Level Matching Algorithm for Clos Network (d-MAC) / 395
- 12.9 The ATLANTA Switch / 398
- 12.10 Concurrent Round-Robin Dispatching (CRRD) Scheme / 400
- 12.11 The Path Switch / 404
 - 12.11.1 Homogeneous Capacity and Route Assignment / 406
 - 12.11.2 Heterogeneous Capacity Assignment / 408

13 MULTI-PLANE MULTI-STAGE BUFFERED SWITCH

413

- 13.1 TrueWay Switch Architecture / 414
 - 13.1.1 Stages of the Switch / 415
- 13.2 Packet Scheduling / 417
 - 13.2.1 Partial Packet Interleaving (PPI) / 419
 - 13.2.2 Dynamic Packet Interleaving (DPI) / 419
 - 13.2.3 Head-of-Line (HOL) Blocking / 420
- 13.3 Stage-To-Stage Flow Control / 420
 - 13.3.1 Back-Pressure / 421
 - 13.3.2 Credit-Based Flow Control / 421
 - 13.3.3 The DQ Scheme / 422
- 13.4 Port-To-Port Flow Control / 424
 - 13.4.1 Static Hashing / 424
 - 13.4.2 Dynamic Hashing / 425
 - 13.4.3 Time-Stamp-Based Resequence / 428
 - 13.4.4 Window-Based Resequence / 428

- 13.5 Performance Analysis / 431
 - 13.5.1 Random Uniform Traffic / 431
 - 13.5.2 Hot-Spot Traffic / 432
 - 13.5.3 Bursty Traffic / 432
 - 13.5.4 Hashing Schemes / 432
 - 13.5.5 Window-Based Resequencing Scheme / 434
- 13.6 Prototype / 434

14 LOAD-BALANCED SWITCHES 438

- 14.1 Birkhoff–Von Neumann Switch / 438
- 14.2 Load-Balanced Birkhoff–von Neumann Switches / 441
 - 14.2.1 Load-Balanced Birkhoff–von Neumann Switch Architecture / 441
 - 14.2.2 Performance of Load-Balanced Birkhoff–von Neumann Switches / 442
- 14.3 Load-Balanced Birkhoff–von Neumann Switches With FIFO Service / 444
 - 14.3.1 First Come First Served (FCFS) / 446
 - 14.3.2 Earliest Deadline First (EDF) and EDF-3DQ / 450
 - 14.3.3 Full Frames First (FFF) / 451
 - 14.3.4 Full Ordered Frames First (FOFF) / 455
 - 14.3.5 Mailbox Switch / 456
 - 14.3.6 Byte-Focal Switch / 459

15 OPTICAL PACKET SWITCHES 468

- 15.1 Opto-Electronic Packet Switches / 469
 - 15.1.1 Hypass / 469
 - 15.1.2 Star-Track / 471
 - 15.1.3 Cisneros and Brackett / 472
 - 15.1.4 BNR (Bell-North Research) Switch / 473
 - 15.1.5 Wave-Mux Switch / 474
- 15.2 Optoelectronic Packet Switch Case Study I / 475
 - 15.2.1 Speedup / 476
 - 15.2.2 Data Packet Flow / 477
 - 15.2.3 Optical Interconnection Network (OIN) / 477
 - 15.2.4 Ping-Pong Arbitration Unit / 482
- 15.3 Optoelectronic Packet Switch Case Study II / 490
 - 15.3.1 Petabit Photonic Packet Switch Architecture / 490
 - 15.3.2 Photonic Switch Fabric (PSF) / 495
- 15.4 All Optical Packet Switches / 503
 - 15.4.1 The Staggering Switch / 503
 - 15.4.2 ATMOS / 504

- 15.4.3 Duan's Switch / 505
- 15.4.4 3M Switch / 506
- 15.5 Optical Packet Switch with Shared Fiber Delay Lines
Single-stage Case / 509
 - 15.5.1 Optical Cell Switch Architecture / 509
 - 15.5.2 Sequential FDL Assignment (SEFA) Algorithm / 512
 - 15.5.3 Multi-Cell FDL Assignment (MUFA) Algorithm / 518
- 15.6 All Optical Packet Switch with Shared Fiber Delay
Lines – Three Stage Case / 524
 - 15.6.1 Sequential FDL Assignment for
Three-Stage OCNS (SEFAC) / 526
 - 15.6.2 Multi-Cell FDL Assignment for
Three-Stage OCNS (MUFAC) / 526
 - 15.6.3 FDL Distribution in Three-Stage OCNS / 528
 - 15.6.4 Performance Analysis of SEFAC and MUFAC / 530
 - 15.6.5 Complexity Analysis of SEFAC and MUFAC / 532

16 HIGH-SPEED ROUTER CHIP SET

538

- 16.1 Network Processors (NPs) / 538
 - 16.1.1 Overview / 538
 - 16.1.2 Design Issues for Network Processors / 539
 - 16.1.3 Architecture of Network Processors / 542
 - 16.1.4 Examples of Network Processors – Dedicated Approach / 543
- 16.2 Co-Processors for Packet Classification / 554
 - 16.2.1 LA-1 Bus / 554
 - 16.2.2 TCAM-Based Classification Co-Processor / 556
 - 16.2.3 Algorithm-Based Classification Co-Processor / 562
- 16.3 Traffic Management Chips / 567
 - 16.3.1 Overview / 567
 - 16.3.2 Agere's TM Chip Set / 567
 - 16.3.3 IDT TM Chip Set / 573
 - 16.3.4 Summary / 579
- 16.4 Switching Fabric Chips / 579
 - 16.4.1 Overview / 579
 - 16.4.2 Switch Fabric Chip Set from Vitesse / 580
 - 16.4.3 Switch Fabric Chip Set from AMCC / 589
 - 16.4.4 Switch Fabric Chip Set from IBM (now of AMCC) / 593
 - 16.4.5 Switch Fabric Chip Set from Agere / 597

INDEX

606

PREFACE

As increasing voice, audio, video, TV, and gaming traffic is carried over IP, Internet traffic continues to grow rapidly. Many network-related applications are emerging for portable devices. As smart cellular phone technology advances, the price decreases, and the infrastructure to support wireless applications (voice, data, video) is being deployed ubiquitously to meet unprecedented demands from users. All of these fast-growing services translate into the high volume of Internet traffic, stringent quality of service (QoS) requirements, large number of hosts/devices to be supported, large forwarding tables to support, high speed packet processing, and large storage capability. When designing/operating next generation switches and routers, these factors create new specifications and new challenges for equipment vendors and network providers.

Jonathan has co-authored two books: *Broadband Packet Switching Technologies—A Practical Guide to ATM Switches and IP Routers* and *Quality of Service Control in High-Speed Networks*, published by John Wiley in 2001. Because the technologies in both electronics and optics have significantly advanced and because the design specifications for routers have become more demanding and challenging, it is time to write another book. This book includes new architectures, algorithms, and implementations developed since 2001. Thus, it is more updated and more complete than the two previous books.

In addition to the need for high-speed and high-capacity transmission/switching equipment, the control function of the equipment and network has also become more sophisticated in order to support new features and requirements of the Internet, including fast re-routing due to link failure (one or more failures), network security, network measurement for dynamic routing, and easy management. This book focuses on the subsystems and devices on the data plane. There is a brief introduction to IP network management to familiarize readers with how the network is managed, as many routers are interconnected together.

The book starts with an introduction to today's and tomorrow's networks, the router architectures and their building blocks, examples of commercial high-end routers, and the challenging issues of designing high-performance high-speed routers. The book first covers the main functions in the line cards of a core router, including route lookup, packet classification, and traffic management for QoS control described in Chapters 2, 3, and

4, respectively. It then follows with 11 chapters in packet switching designs, covering various architectures, algorithms, and technologies (including electrical and optical packet switching). The last chapter of the book presents the state-of-the-art commercial chipsets used to build the routers. This is one of the important features in this book—showing readers the architecture and functions of practical chipsets to reinforce the theories and conceptual designs covered in previous chapters.

A distinction of this book is that we provide as many figures as possible to explain the concepts. Readers are encouraged to first scan through the figures and try to understand them before reading the text. If fully understood, readers can skip to the text to save time. However, the text is written in such a way as to talk the readers through the figures.

Jonathan and Bin each have about 20 years of experience researching high-performance switches and routers, implementing them in various systems with VLSI (very-large-scale integration) and FPGA (field-programmable gate array) chips, transferring technology to the industry, and teaching such subjects in the college and to the industry companies. They have accumulated their practical experience in writing this book. The book includes theoretical concepts and algorithms, design architectures, and actual implementations. It will benefit the readers in different aspects of building a high-performance switch/router. The draft of the book has been used as a text for the past two years when teaching senior undergraduate and first-year graduate students at the author's universities. If any errors are found, please send an email to chao@poly.edu. The authors will then make the corresponding corrections in future editions.

Audience

This book is an appropriate text for senior and graduate students in Electrical Engineering, Computer Engineering, and Computer Science. They can embrace the technology of the Internet so as to better position themselves when they graduate and look for jobs in the high-speed networking field. This book can also be used as a reference for people working in the Internet-related area. Engineers from network equipment vendors and service providers can also benefit from the book by understanding the key concepts of packet switching systems and the key techniques of building high-speed and high-performance routers.

ACKNOWLEDGMENTS

This book would not have been published without the help of many people. We would like to thank them for their efforts in improving the quality of the book.

Several chapters of the book are based on research work that was done at Polytechnic University and Tsinghua University. We would like to thank several individuals who contributed material to some sections. They are Professor Ming Yu (Florida State University) on Section 1.5, Professor Derek C. W. Pao (City University of Hong Kong) on Section 2.4.2, and Professor Aleksandra Smiljanic (Belgrade University) on a scheduling scheme she proposed in Chapter 7. We would like to express our gratitude to Dr. Yihan Li (Auburn University) for her contribution to part of Chapter 7, and the students in Bin's research group in Tsinghua University for their contribution to some chapters. They are Chenchen Hu, Kai Zheng, Zhen Liu, Lei Shi, Xuefei Chen, Xin Zhang, Yang Xu, Wenjie Li, and Wei Li. The manuscript has been managed from the beginning to the end by Mr Jian Li (Polytechnic University), who has put in tremendous effort to carefully edit the manuscript and serve as a coordinator with the publisher.

The manuscript draft was reviewed by the following people and we would like to thank them for their valuable feedback: Professor Cristina López Bravo (University of Vigo, Spain), Dr Hiroaki Harai (Institute of Information and Communications Technology, Japan), Dr Simin He (Chinese Academy of Sciences), Professor Hao Che (University of Texas at Arlington), Professor Xiaohong Jiang (Tohoku University, Japan), Dr Yihan Li (Auburn University), Professor Dr Soung Yue Liew (Universiti Tunku Abdul Rahman, Malaysia), Dr Jan van Lunteren (IBM, Zurich), Professor Jinsoo Park (Essex County College, New Jersey), Professor Roberto Rojas-cessa (New Jersey Institute of Technology), Professor Aleksandra Smiljanic (Belgrade University, Serbia and Montenegro), Professor Dapeng Wu (University of Florida), and Professor Naoaki Yamanaka (Keio University, Japan).

Jonathan would like to thank his wife, Ammie, and his children, Jessica, Roger, and Joshua, for their love, support, encouragement, patience, and perseverance. He also thanks his parents for their encouragement.

Bin would like to thank his wife, Yingjun Ma, and his daughter, Jenny for their understanding and support. He also thanks his father-in-law for looking after Jenny to spare his time to prepare the book.

CHAPTER 1

INTRODUCTION

The Internet, with its robust and reliable Internet Protocol (IP), is widely considered the most reachable platform for the current and next generation information infrastructure. The virtually unlimited bandwidth of optical fiber has tremendously increased the data transmission speed over the past decade. Availability of unlimited bandwidth has stimulated high-demand multimedia services such as distance learning, music and video download, and videoconferencing. Current broadband access technologies, such as digital subscriber lines (DSLs) and cable television (CATV), are providing affordable broadband connection solutions to the Internet from home. Furthermore, with Gigabit Ethernet access over dark fiber to the enterprise on its way, access speeds are expected to largely increase. It is clear that the deployment of these broadband access technologies will result in a high demand for large Internet bandwidth. To keep pace with the Internet traffic growth, researchers are continually exploring faster transmission and switching technologies. The advent of optical transmission technologies, such as dense wave division multiplexing (DWDM), optical add-drop multiplexers, and ultra-long-haul lasers have had a large influence on lowering the costs of digital transmission. For instance, 300 channels of 11.6 Gbps can be wavelength-division multiplexed on a single fiber and transmitted over 7000 km [1]. In addition, a 1296×1296 optical cross-connect (OXC) switching system using micro-electro-mechanical systems (MEMS) with a total switching capacity of 2.07 petabits/s has been demonstrated [2]. In the rest of this chapter, we explore state-of-the-art network infrastructure, future design trends, and their impact on next generation routers. We also describe router architectures and the challenges involved in designing high-performance large-scale routers.

1.1 ARCHITECTURE OF THE INTERNET: PRESENT AND FUTURE

1.1.1 The Present

Today's Internet is an amalgamation of thousands of commercial and service provider networks. It is not feasible for a single service provider to connect two distant nodes on the Internet. Therefore, service providers often rely on each other to connect the dots. Depending on the size of network they operate, Internet Service Providers (ISPs) can be broken down into three major categories. Tier-1 ISPs are about a dozen major telecommunication companies, such as UUNet, Sprint, Qwest, XO Network, and AT&T, whose high-speed global networks form the Internet backbone. Tier-1 ISPs do not buy network capacity from other providers; instead, they sell or lease access to their backbone resource to smaller Tier-2 ISPs, such as America Online and Broadwing. Tier-3 ISPs are typically regional service providers such as Verizon and RCN through whom most enterprises connect to the Internet. Figure 1.1 illustrates the architecture of a typical Tier-1 ISP network.

Each Tier-1 ISP operates multiple IP/MPLS (multi-protocol label switching), and sometimes ATM (asynchronous transfer mode), backbones with speeds varying anywhere from T3 to OC-192 (optical carrier level 192, ~10 Gbps). These backbones are interconnected through peering agreements between ISPs to form the Internet backbone. The backbone is designed to transfer large volumes of traffic as quickly as possible between networks. Enterprise networks are often linked to the rest of the Internet via a variety of links, anywhere from a T1 to multiple OC-3 lines, using a variety of Layer 2 protocols, such as Gigabit Ethernet, frame relay, and so on. These enterprise networks are then overhauled into service provider networks through edge routers. An edge router can aggregate links from multiple enterprises. Edge routers are interconnected in a pool, usually at a Point of Presence (POP)

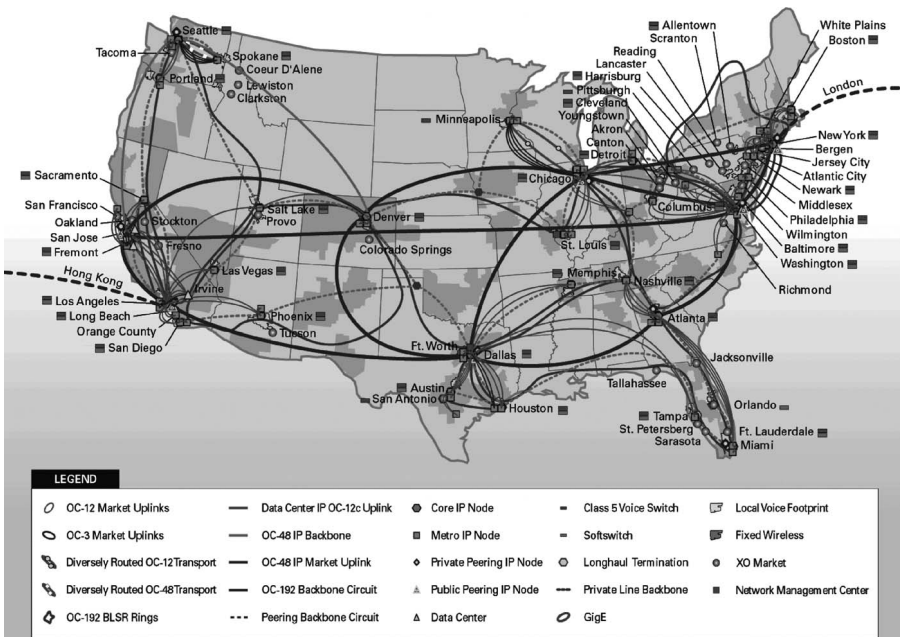


Figure 1.1 Network map of a Tier-1 ISP, XO Network.

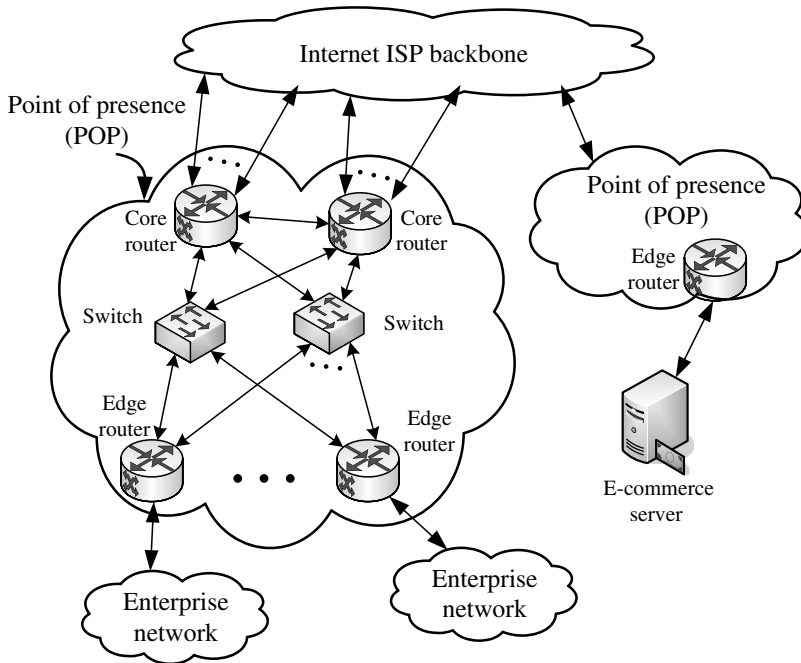


Figure 1.2 Point of presence (POP).

of a service provider, as shown in Figure 1.2. Each POP may link to other POPs of the same ISP through optical transmission/switching equipment, may link to POPs of other ISPs to form a peering, or link to one or more backbone routers. Typically, a POP may have a few backbone routers in a densely connected mesh. In most POPs, each edge router connects to at least two backbone routers for redundancy. These backbone routers may also connect to backbone routers at other POPs according to ISP peering agreements. Peering occurs when ISPs exchange traffic bound for each other's network over a direct link without any fees. Therefore, peering works best when peers exchange roughly the same amount of traffic. Since smaller ISPs do not have high quantities of traffic, they often have to buy *transit* from a Tier-1 provider to connect to the Internet. A recent study of the topologies of 10 service providers across the world shows that POPs share this generic structure [3].

Unlike POPs, the design of backbone varies from service provider to service provider. For example, Figure 1.3 illustrates backbone design paradigms of three major service providers

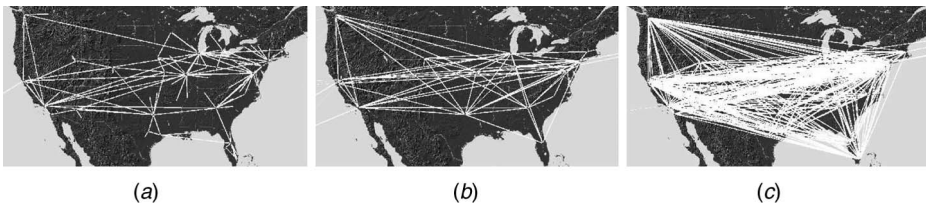


Figure 1.3 Three distinct backbone design paradigms of Tier-1 ISPs. (a) AT&T; (b) Sprint; (c) Level 3 national network infrastructure [3].

in the US. AT&T's backbone design includes large POPs at major cities, which in turn fan out into smaller per-city POPs. In contrast, Sprint's backbone has only 20 well connected POPs in major cities and suburban links are back-hauled into the POPs via smaller ISPs. Most major service providers still have the AT&T backbone model and are in various stages of moving to Sprint's design. Sprint's backbone design provides a good solution to service providers grappling with a need to reduce capital expenditure and operational costs associated with maintaining and upgrading network infrastructure. Interestingly, Level 3 presents another design paradigm in which the backbone is highly connected via circuit technology such as, MPLS, ATM or frame relays. As will be seen later, this is the next generation of network design where the line between backbone and network edge begins to blur.

Now, let us see how network design impacts on the next generation routers. Router design is often guided by the economic requirements of service providers. Service providers would like to reduce the infrastructure and maintenance costs while, at the same time, increasing available bandwidth and reliability. To this end, network backbone has a set of well-defined, narrow requirements. Routers in the backbone should simply move traffic as fast as possible. Network edge, however, has broad and evolving requirements due simply to the diversity of services and Layer 2 protocols supported at the edge. Today most POPs have multiple edge routers optimized for point solutions. In addition to increasing infrastructure and maintenance costs, this design also increases the complexity of POPs resulting in an unreliable network infrastructure. Therefore, newer edge routers have been designed to support diversity and are easily adaptable to the evolving requirements of service providers. This design trend is shown in Table 1.1, which lists some properties of enterprise, edge, and core routers currently on the market. As we will see in the following sections, future network designs call for the removal of edge routers altogether and their replacement with fewer core routers to increase reliability, throughput, and to reduce costs. This means next generation routers would have to amalgamate the diverse service requirements of edge routers and the strict performance requirements of core routers, seamlessly into one body. Therefore, the real question is not whether we should build highly-flexible, scalable, high-performance routers, but how?

1.1.2 The Future

As prices of optical transport and optical switching sharply decrease, some network designers believe that the future network will consist of many mid-size IP routers or MPLS

TABLE 1.1 Popular Enterprise, Edge, and Core Routers in the Market

Model	Capacity ^a	Memory	Power	Features
Cisco 7200	–	256 MB	370 W	QoS, MPLS, Aggregation
Cisco 7600	720 Gbps	1 GB	–	QoS, MPLS, Shaping
Cisco 10000	51.2 Gbps	–	1200 W	QoS, MPLS
Cisco 12000	1.28 Tbps	4 GB	4706 W	MPLS, Peering
Juniper M-320	320 Gbps	2 GB	3150 W	MPLS, QoS, VPN
Cisco CRS	92 Tbps	4 GB	16,560 W	MPLS, Qos, Peering
Juniper TX/T-640	2.5 Tbps/640 Gbps	2 GB	4550 W/6500 W	MPLS, QoS, Peering

^aNote that the listed capacity is the combination of ingress and egress capacities.

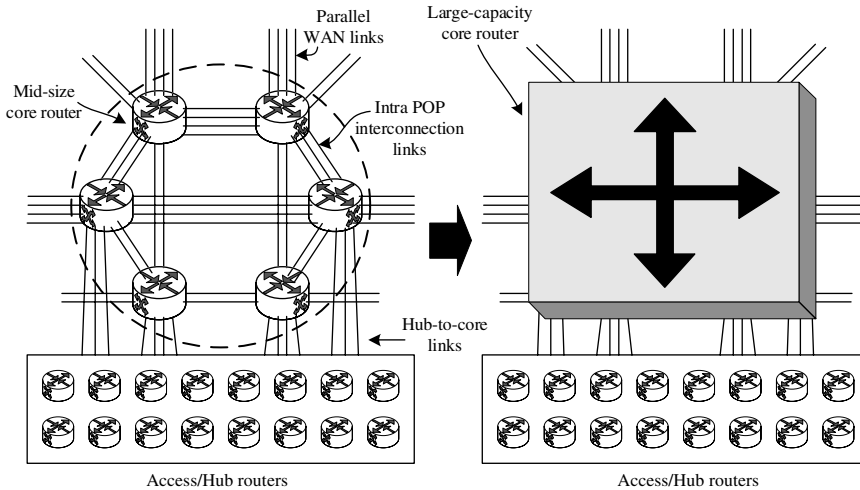


Figure 1.4 Replacing a cluster of mid-size routers with a large-capacity scalable router.

switches at the network edge that are connected to optical crossconnects (OXC), which are then interconnected by DWDM transmission equipment. The problem for this approach is that connections to the OXC are usually high bit rates, for example, 10 Gbps for now and 40 Gbps in the near future. When the edge routers want to communicate with all other routers, they either need to have direct connections to those routers or connect through multiple logical hops (i.e., routed by other routers). The former case results in low link utilization while the latter results in higher latency. Therefore, some network designers believe it is better to build very large IP routers or MPLS switches at POPs. They aggregate traffic from edge routers onto high-speed links that are then directly connected to other large routers at different POPs through DWDM transmission equipment. This approach achieves higher link utilization and fewer hops (thus lower latency). As a result, the need for an OXC is mainly for provisioning and restoring purposes but not for dynamic switching to achieve higher link utilization.

Current router technologies available in the market cannot provide large switching capacities to satisfy current and future bandwidth demands. As a result, a number of mid-size core routers are interconnected with numerous links and use many expensive line cards that are used to carry intra-cluster traffic rather than revenue-generating users' or wide-area-network (WAN) traffic. Figure 1.4 shows how a router cluster is replaced by a large-capacity scalable router, saving the cost of numerous line cards and links, and real estate. It provides a cost-effective solution that can satisfy Internet traffic growth without having to replace routers every two to three years. Furthermore, there are fewer individual routers that need to be configured and managed, resulting in a more efficient and reliable system.

1.2 ROUTER ARCHITECTURES

IP routers' functions can be classified into two categories: datapath functions and control plane functions [4].

The datapath functions such as forwarding decision, forwarding through the backplane, and output link scheduling are performed on every datagram that passes through the router. When a packet arrives at the forwarding engine, its destination IP address is first masked by the subnet mask (logical AND operation) and the resulting address is used to lookup the forwarding table. A so-called longest prefix matching method is used to find the output port. In some applications, packets are classified based on 104 bits that include the IP source/destination addresses, transport layer port numbers (source and destination), and type of protocol, which is generally called 5-tuple. Based on the result of classification, packets may be either discarded (firewall application) or handled at different priority levels. Then, time-to-live (TTL) value is decremented and a new header checksum is recalculated.

The control plane functions include the system configuration, management, and exchange of routing table information. These are performed relatively infrequently. The route controller exchanges the topology information with other routers and constructs a routing table based on a routing protocol, for example, RIP (Routing Information Protocol), OSPF (Open Shortest Path Forwarding), or BGP (Border Gateway Protocol). It can also create a forwarding table for the forwarding engine. Since the control functions are not performed on each arriving individual packet, they do not have a strict speed constraint and are implemented in software in general.

Router architectures generally fall into two categories: centralized (Fig. 1.5a) and distributed (Fig. 1.5b).

Figure 1.5a shows a number of network interfaces, forwarding engines, a route controller (RC), and a management controller (MC) interconnected by a switch fabric. Input interfaces send packet headers to the forwarding engines through the switch fabric. The forwarding engines, in turn, determine which output interface the packet should be sent to. This information is sent back to the corresponding input interface, which forwards the packet to the right output interface. The only task of a forwarding engine is to process packet headers and is shared by all the interfaces. All other tasks such as participating in routing protocols, reserving resource, handling packets that need extra attention, and other administrative and maintenance tasks, are handled by the RC and the MC. The BBN multi-gigabit router [5] is an example of this design.

The difference between Figure 1.5a and 1.5b is that the functions of the forwarding engines are integrated into the interface cards themselves. Most high-performance routers use this architecture. The RC maintains a routing table and updates it based on routing protocols used. The routing table is used to generate a forwarding table that is then downloaded

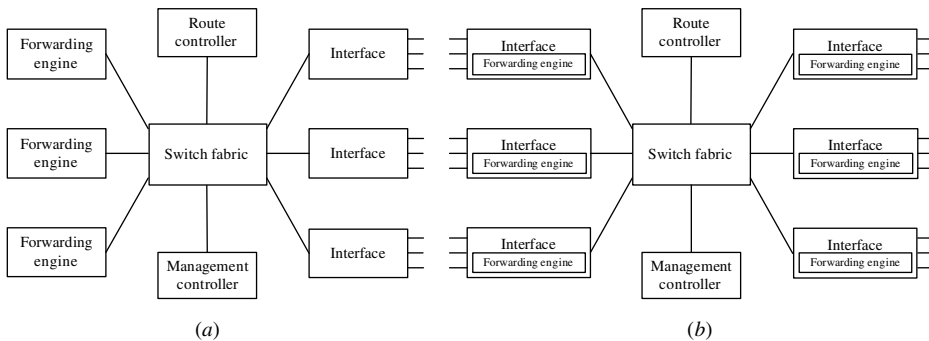


Figure 1.5 (a) Centralized versus (b) distributed models for a router.

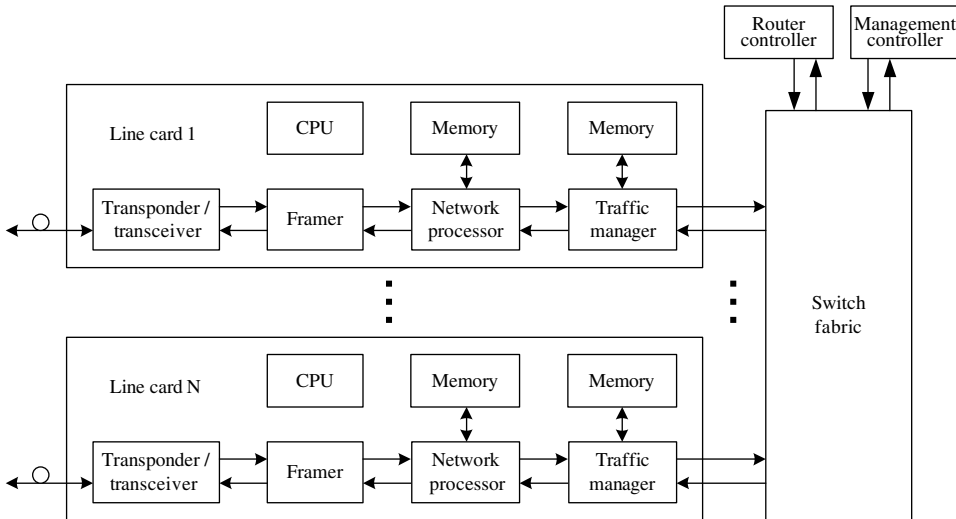


Figure 1.6 Typical router architecture.

from the RC to the forwarding engines in the interface cards. It is not necessary to download a new forwarding table for every route update. Route updates can be frequent, but routing protocols need time, in the order of minutes, to converge. The RC needs a dynamic routing table designed for fast updates and fast generation of forwarding tables. Forwarding tables, on the other hand, can be optimized for lookup speed and need not be dynamic.

Figure 1.6 shows a typical router architecture, where multiple line cards, an RC, and an MC are interconnected through a switch fabric. The communication between the RC/MC and the line cards can be either through the switch fabric or through a separate interconnection network, such as a Ethernet switch. The line cards are the entry and exit points of data to and from a router. They provide the interface from physical and higher layers to the switch fabric. The tasks provided by line cards are becoming more complex as new applications develop and protocols evolve. Each line card supports at least one full-duplex fiber connection on the network side, and at least one ingress and one egress connection to the switch fabric backplane. Generally speaking, for high-bandwidth applications, such as OC-48 and above, the network connections support channelization for aggregation of lower-speed lines into a large pipe, and the switch fabric connections provide flow-control mechanisms for several thousand input and output queues to regulate the ingress and egress traffic to and from the switch fabric.

A line card usually includes components such as a transponder, framer, network processor (NP), traffic manager (TM), and central processing unit (CPU).

Transponder/Transceiver: This component performs optical-to-electrical and electrical-to-optical signal conversions, and serial-to-parallel and parallel-to-serial conversions [6, 7].

Framer: A framer performs synchronization, frame overhead processing, and cell or packet delineation. On the transmit side, a SONET (synchronous optical network)/SDH (synchronous digital hierarchy) framer generates section, line, and path overhead. It performs framing pattern insertion (A1, A2) and scrambling. It

generates section, line, and path bit interleaved parity (B1/B2/B3) for far-end performance monitoring. On the receive side, it processes section, line, and path overhead. It performs frame delineation, descrambling, alarm detection, pointer interpretation, bit interleaved parity monitoring (B1/B2/B3), and error count accumulation for performance monitoring [8]. An alternative for the framer is Ethernet framer.

Network Processor: The NP mainly performs table lookup, packet classification, and packet modification. Various algorithms to implement the first two functions are presented in Chapters 2 and 3, respectively. The NP can perform those two functions at the line rate using external memory, such as static random access memory (SRAM) or dynamic random access memory (DRAM), but it may also require external content addressable memory (CAM) or specialized co-processors to perform deep packet classification at higher levels. In Chapter 16, we present some commercially available NP and ternary content addressable memory (TCAM) chips.

Traffic Manager: To meet the requirements of each connection and service class, the TM performs various control functions to cell/packet streams, including traffic access control, buffer management, and cell/packet scheduling. Traffic access control consists of a collection of specification techniques and mechanisms that (1) specify the expected traffic characteristics and service requirements (e.g., peak rate, required delay bound, loss tolerance) of a data stream; (2) shape (i.e., delay) data streams (e.g., reducing their rates and/or burstiness); and (3) police data streams and take corrective actions (e.g., discard, delay, or mark packets) when traffic deviates from its specification. The usage parameter control (UPC) in ATM and differentiated service (DiffServ) in IP performs similar access control functions at the network edge. Buffer management performs cell/packet discarding, according to loss requirements and priority levels, when the buffer exceeds a certain threshold. Proposed schemes include early packet discard (EPD) [9], random early packet discard (REPD) [10], weighted REPD [11], and partial packet discard (PPD) [12]. Packet scheduling ensures that packets are transmitted to meet each connection's allocated bandwidth/delay requirements. Proposed schemes include deficit round-robin, weighted fair queuing (WFQ) and its variants, such as shaped virtual clock [13] and worst-case fairness WFQ (WF^2Q+) [14]. The last two algorithms achieve the worst-case fairness properties. Details are discussed in Chapter 4. Many quality of service (QoS) control techniques, algorithms, and implementation architectures can be found in Ref. [15]. The TM may also manage many queues to resolve contention among the inputs of a switch fabric, for example, hundreds or thousands of virtual output queues (VOQs). Some of the representative TM chips on the market are introduced in Chapter 16, whose purpose it is to match the theories in Chapter 4 with practice.

Central Processing Unit. The CPU performs control plane functions including connection set-up/tear-down, table updates, register/buffer management, and exception handling. The CPU is usually not in-line with the fast-path on which maximum-bandwidth network traffic moves between the interfaces and the switch fabric.

The architecture in Figure 1.6 can be realized in a multi-rack (also known as multi-chassis or multi-shelf) system as shown in Figure 1.7. In this example, a half rack, equipped with a switch fabric, a duplicated RC, a duplicated MC, a duplicated system clock (CLK), and a duplicated fabric shelf controller (FSC), is connected to all other line card (LC) shelves, each of which has a duplicated line card shelf controller (LSC). Both the FSC and the

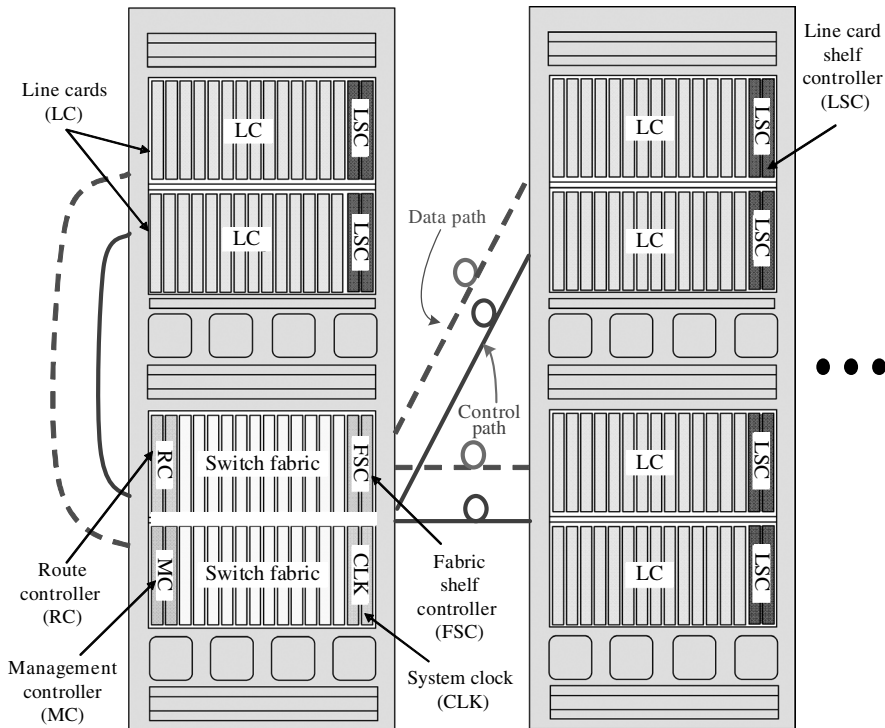


Figure 1.7 Multi-rack router system.

LSC provide local operation and maintenance for the switch fabric and line card shelves, respectively. They also provide the communication channels between the switch/line cards with the RC and the MC. The duplicated cards are for reliability concerns. The figure also shows how the system can grow by adding more LC shelves. Interconnections between the racks are sets of cables or fibers, carrying information for the data and the control planes. The cabling usually is a combination of unshielded twisted path (UTP) Category 5 Ethernet cables for control path, and fiber-optic arrays for data path.

1.3 COMMERCIAL CORE ROUTER EXAMPLES

We now briefly discuss the two most popular core routers on the market: Juniper Network's T640 TX-Matrix [16] and Cisco System's Carrier Routing System (CRS-1) [17].

1.3.1 T640 TX-Matrix

A T640 TX-Matrix is composed of up to four routing nodes and a TX Routing Matrix interconnecting the nodes. A TX Routing Matrix connects up to four T640 routing nodes via a three-stage Clos network switch fabric to form a unified router with the capacity of 2.56 Terabits. The blueprint of a TX Routing Matrix is shown in Figure 1.8. The unified router is controlled by the Routing Engine of the matrix which is responsible for running routing protocols and for maintaining overall system state. Routing engines in each routing

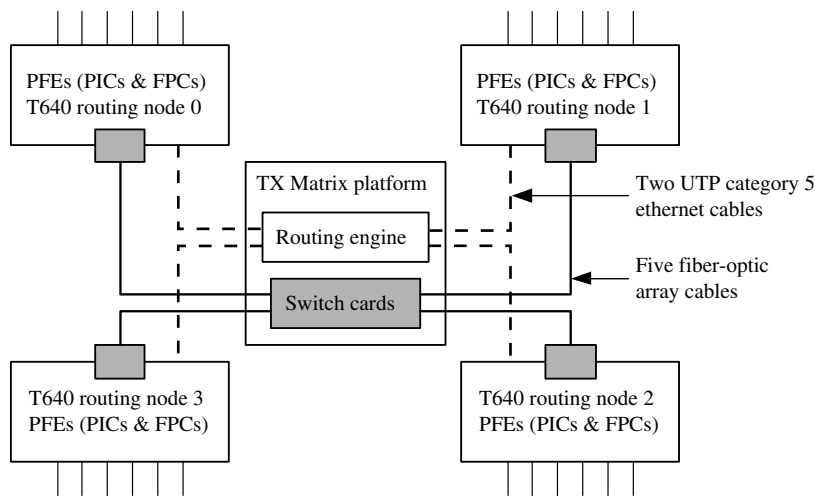


Figure 1.8 TX Routing Matrix with four T640 routing nodes.

node manage their individual components in coordination with the routing engine of the matrix. Data and control plane of each routing node is interconnected via an array of optical and Ethernet cables. Data planes are interconnected using VCSEL (vertical cavity surface emitting laser) optical lines whereas control planes are interconnected using UTP Category 5 Ethernet cables.

As shown in Figure 1.9, each routing node has two fundamental architectural components, namely the control plane and the data plane. The T640 routing node’s control plane is implemented by the JUNOS software that runs on the node’s routing engine. JUNOS is a micro-kernel-based modular software that assures reliability, fault isolation, and high availability. It implements the routing protocols, generates routing tables and forwarding tables, and supports the user interface to the router. Data plane, on the other hand, is responsible for processing packets in hardware before forwarding them across the switch fabric from the ingress interface to the appropriate egress interface. The T640 routing node’s data plane is implemented in custom ASICs in a distributed architecture.

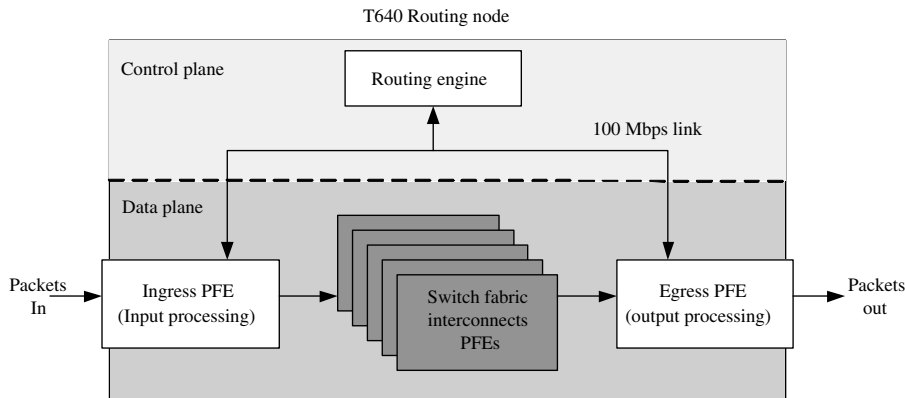


Figure 1.9 T640 routing node architecture.