# Network Systems Design

using

# Network Processors



DOUGLAS E. COMER

# Contents

## Chapter 3   Review Of Protocols And Packet Formats                    15

# PART I   Traditional Protocol Processing Systems

## Chapter 4   Conventional Computer Hardware Architecture              29

## Chapter 5   Basic Packet Processing: Algorithms And Data Structures   43

## Chapter 6    Packet Processing Functions                                             67

## Chapter 7    Protocol Software On A Conventional Processor                   83

## Chapter 8  Hardware Architectures For Protocol Processing                    97

## Chapter 9  Classification And Forwarding                                      115

## Chapter 10  Switching Fabrics                                                        **133**

# PART II  Network Processor Technology

## Chapter 11  Network Processors: Motivation And Purpose                 **153**

## Chapter 12  The Complexity Of Network Processor Design                165

## Chapter 13  Network Processor Architectures                177

## Chapter 14  Issues In Scaling A Network Processor                195

## Chapter 15  Examples Of Commercial Network Processors          213

## Chapter 16  Languages Used For Classification          233

## Chapter 17  Design Tradeoffs And Consequences                    261

# PART III  Example Network Processor


## Chapter 18  Overview Of The Intel Network Processor                273

# Chapter 23  ACE Run-Time Structure And StrongARM Facilities 349

# Chapter 24  Microengine Programming I 371

## Chapter 27  Intel's Second Generation Processors                          **441**

## Appendix 1  Glossary Of Terms And Abbreviations                          **449**

## Bibliography                                                                **497**

## Index                                                                       **501**