

Quantitative System Performance

COMPUTER
SYSTEM
ANALYSIS
USING
QUEUEING
NETWORK
MODELS

Edward D. Lazowska | John Zahorjan
G. Scott Graham | Kenneth C. Sevcik

Quantitative System Performance

Computer System Analysis Using Queueing Network Models

Edward D. Lazowska, John Zahorjan,
G. Scott Graham, and Kenneth C. Sevcik

Prentice-Hall, Inc., Englewood Cliffs, New Jersey 07632

Library of Congress Cataloging in Publication Data

Main entry under title:

Quantitative system performance.

Bibliography: p.

Includes index.

1. Electronic digital computers—Evaluation.
2. Digital computer simulation. 3. Queuing theory.

I. Lazowska, Edward D. (date)

QA76.9.E94Q36 1984 001.64'028'7 83-13791

ISBN 0-13-746975-6

© 1984 by Prentice-Hall, Inc., Englewood Cliffs, New Jersey 07632

All rights reserved. No part of this book may be reproduced in any form or by any means without permission in writing from the publisher.

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

ISBN 0-13-746975-6

Editorial/production supervision: Nancy Milnamow

Cover design: Edsal Enterprises

Jacket design: Lundgren Graphics, Ltd.

Manufacturing buyer: Gordon Osbourne

Prentice-Hall International, Inc., *London*
Prentice-Hall of Australia Pty. Limited, *Sydney*
Editora Prentice-Hall do Brasil, *Rio de Janeiro*
Prentice-Hall of Canada Inc., *Toronto*
Prentice-Hall of India Private Limited, *New Delhi*
Prentice-Hall of Japan, Inc., *Tokyo*
Prentice-Hall of Southeast Asia Pte. Ltd., *Singapore*
Whitehall Books Limited, *Wellington, New Zealand*

Quantitative System Performance

COMPUTER SYSTEM ANALYSIS USING
QUEUEING NETWORK MODELS

Edward D. Lazowska • John Zahorjan
G. Scott Graham • Kenneth C. Sevcik

Comprehensive coverage of the computer system performance evaluation technology of the 1980s—queueing network models. A practical, accessible book designed to teach professional computer system performance analysts how to apply queueing network models in their work, as tools to assist in answering the questions of cost and performance that arise throughout the life of a computer system: during design and implementation, during sizing and acquisition, and during evolution of the configuration and workload.

Uses the *operational* approach to queueing network modelling, emphasizing practical rather than theoretical considerations. Case studies and examples are integrated throughout. Discusses the structure and use of queueing network modelling software packages. Includes exercises, helpful for self-study or in a classroom setting. Fortran programs implement key algorithms. A self-contained presentation, relying only on experience with computer systems.

(continued from front flap)
Edward D. Lazowska and **John Zahorjan** are with the Department of Computer Science at the University of Washington. **G. Scott Graham** and **Kenneth C. Sevcik** are with the Computer Systems Research Group at the University of Toronto. This book reflects their extensive experience in the theory and practice of computer system analysis using queueing network models:

- the development of commercial and research queueing network modelling software
- consulting in the public and private sectors
- the development and teaching of graduate courses and of professional short courses, seminars, and tutorials
- research, ranging from relatively applied to relatively theoretical
- service, in capacities such as journal editorship and conference organization

PRENTICE-HALL, Inc.
Englewood Cliffs, N.J. 07632

Printed in the U.S. of America

13374
1987

Quantitative System Performance
Computer System Analysis Using
Queueing Network Models

Preface

This book is written for computer system performance analysts. Its goal is to teach them to apply queueing network models in their work, as tools to assist in answering the questions of cost and performance that arise throughout the life of a computer system.

Our approach to the subject arises from our collective experience in contributing to the theory of queueing network modelling, in embodying this theory in performance analysis tools, in applying these tools in the field, and in teaching computer system analysis using queueing network models in academic and industrial settings. Some important beliefs underlying our approach are:

- Although queueing network models are not a panacea, they are the appropriate tool in a wide variety of computer system design and analysis applications.
- The single most important attribute of a computer system analyst is a thorough understanding of computer systems. We assume this of our readers.
- On the one hand, mathematical sophistication is not required to analyze computer systems intelligently and successfully using queueing network models. This is the case because the algorithms for evaluating queueing network models are well developed.
- On the other hand, the purchase of a queueing network modelling software package does not assure success in computer system analysis. This is the case because defining and parameterizing a queueing network model of a particular computer system is a blend of art and science, requiring training and experience.

Queueing network modelling is a methodology for the analysis of computer systems. A methodology is a way of thinking, not a substitute for thinking.

We have divided the book into six parts. In Part I we provide four types of background material: a general discussion of queueing network modelling, an overview of the way in which a modelling study is conducted, an introduction to the interesting performance quantities in computer systems and to certain relationships that must hold among them, and a discussion of the inputs and outputs of queueing network models.

In Part II we present the techniques that are used to evaluate queueing network models — to obtain outputs such as utilizations, residence times, queue lengths, and throughputs from inputs such as workload intensities and service demands.

In Part III we explore the need for detailed models of specific subsystems, and the construction of such models for memory, disk I/O, and processor subsystems.

In Part IV we study the parameterization of queueing network models of existing systems, evolving systems, and proposed systems.

In Part V we survey some non-traditional applications, such as the analysis of computer communication networks and database concurrency control mechanisms. We also examine the structure and use of queueing network modelling software packages.

In Part VI, the appendices, we provide a case study in obtaining queueing network parameter values from system measurement data, and programs implementing the queueing network evaluation techniques described in Part II.

Case studies appear throughout the book. They are included to illustrate various aspects of computer system analysis using queueing network models. They should *not* be misconstrued as making general statements about the relative performance of various systems; the results have significance only for the specific configurations and workloads under consideration.

We have summarized a number of important modelling techniques in the form of “Algorithms”. Our intention is to provide enough information that the reader can understand fully the essential aspects of each technique. We omit details of significance to the implementation of a technique when we feel that these details might obscure the more fundamental concepts.

It is our experience that practicing computer system analysts are relatively skilled in techniques such as workload characterization, system measurement, interpretation of performance data, and system tuning, and are at least acquainted with basic statistical methods and with simulation. Each of these subjects is well represented in the existing literature, and is given short shrift in the present book. Much interesting and important research work concerning queueing network modelling also is given short shrift; we discuss the one approach to each problem that we feel is best suited for application. For readers who desire to pursue a topic in greater detail than we have provided, each chapter concludes with a brief discussion of the relevant literature.

We owe a significant debt to Jeffrey P. Buzen and Peter J. Denning, who have been instrumental in the development of a pragmatic

philosophy of computer system analysis using queueing network models. Their influence is evident especially in our use of the *operational* framework for queueing network modelling, which conveys much greater intuition than the more traditional *stochastic* framework.

Jeffrey A. Brumfield, Jeffrey P. Buzen, Domenico Ferrari, Lenny Freilich, and Roger D. Stoesz have assisted us by reviewing our manuscript, as have several anonymous reviewers. Our work in computer system analysis using queueing network models has been supported in part by the National Science Foundation and by the Natural Sciences and Engineering Research Council of Canada. We thank our colleagues at the University of Washington and at the University of Toronto for their encouragement, and our families and friends for their forbearance.

Edward D. Lazowska, John Zahorjan,
G. Scott Graham, and Kenneth C. Sevcik

Seattle and Toronto

Contents

Preface	xii
I. Preliminaries	1
1. An Overview of Queueing Network Modelling	2
1.1. <i>Introduction ...</i>	2
1.2. <i>What Is a Queueing Network Model? ...</i>	4
1.3. <i>Defining, Parameterizing, and Evaluating Queueing Network Models ...</i>	9
1.4. <i>Why Are Queueing Network Models Appropriate Tools? ...</i>	14
1.5. <i>Related Techniques ...</i>	14
1.6. <i>Summary ...</i>	17
1.7. <i>References ...</i>	17
2. Conducting a Modelling Study	20
2.1. <i>Introduction ...</i>	20
2.2. <i>The Modelling Cycle ...</i>	22
2.3. <i>Understanding the Objectives of a Study ...</i>	27
2.4. <i>Workload Characterization ...</i>	30
2.5. <i>Sensitivity Analysis ...</i>	33
2.6. <i>Sources of Insight ...</i>	35
2.7. <i>Summary ...</i>	37
2.8. <i>References ...</i>	39
3. Fundamental Laws	40
3.1. <i>Introduction ...</i>	40
3.2. <i>Basic Quantities ...</i>	40
3.3. <i>Little's Law ...</i>	42
3.4. <i>The Forced Flow Law ...</i>	47
3.5. <i>The Flow Balance Assumption ...</i>	51
3.6. <i>Summary ...</i>	52
3.7. <i>References ...</i>	53
3.8. <i>Exercises ...</i>	54

4.	Queueing Network Model Inputs and Outputs	57
4.1.	<i>Introduction ...</i>	57
4.2.	<i>Model Inputs ...</i>	57
4.3.	<i>Model Outputs ...</i>	60
4.4.	<i>Multiple Class Models ...</i>	62
4.5.	<i>Discussion ...</i>	64
4.6.	<i>Summary ...</i>	67
4.7.	<i>Exercises ...</i>	68
II.	General Analytic Techniques	69
5.	Bounds on Performance	70
5.1.	<i>Introduction ...</i>	70
5.2.	<i>Asymptotic Bounds ...</i>	72
5.3.	<i>Using Asymptotic Bounds ...</i>	77
5.4.	<i>Balanced System Bounds ...</i>	86
5.5.	<i>Summary ...</i>	92
5.6.	<i>References ...</i>	94
5.7.	<i>Exercises ...</i>	95
6.	Models with One Job Class	98
6.1.	<i>Introduction ...</i>	98
6.2.	<i>Workload Representation ...</i>	99
6.3.	<i>Case Studies ...</i>	102
6.4.	<i>Solution Techniques ...</i>	108
6.5.	<i>Theoretical Foundations ...</i>	119
6.6.	<i>Summary ...</i>	121
6.7.	<i>References ...</i>	123
6.8.	<i>Exercises ...</i>	124
7.	Models with Multiple Job Classes	127
7.1.	<i>Introduction ...</i>	127
7.2.	<i>Workload Representation ...</i>	128
7.3.	<i>Case Studies ...</i>	129
7.4.	<i>Solution Techniques ...</i>	134
7.5.	<i>Theoretical Foundations ...</i>	147
7.6.	<i>Summary ...</i>	149
7.7.	<i>References ...</i>	150
7.8.	<i>Exercises ...</i>	150

8.	Flow Equivalence and Hierarchical Modelling	152
8.1.	<i>Introduction ...</i>	152
8.2.	<i>Creating Flow Equivalent Service Centers ...</i>	155
8.3.	<i>Obtaining the Parameters ...</i>	158
8.4.	<i>Solving the High-Level Models ...</i>	159
8.5.	<i>An Application of Hierarchical Modelling ...</i>	160
8.6.	<i>Summary ...</i>	173
8.7.	<i>References ...</i>	174
8.8.	<i>Exercises ...</i>	175
III.	Representing Specific Subsystems	177
9.	Memory	179
9.1.	<i>Introduction ...</i>	179
9.2.	<i>Systems with Known Average Multiprogramming Level ...</i>	181
9.3.	<i>Memory Constraints ...</i>	184
9.4.	<i>Swapping ...</i>	196
9.5.	<i>Paging ...</i>	201
9.6.	<i>Case Studies ...</i>	206
9.7.	<i>Summary ...</i>	217
9.8.	<i>References ...</i>	218
9.9.	<i>Exercises ...</i>	219
10.	Disk I/O	222
10.1.	<i>Introduction ...</i>	222
10.2.	<i>Channel Contention in Non-RPS I/O Subsystems ...</i>	225
10.3.	<i>Channel Contention in RPS I/O Subsystems ...</i>	230
10.4.	<i>Additional Path Elements ...</i>	233
10.5.	<i>Multipathing ...</i>	237
10.6.	<i>Other Architectural Characteristics ...</i>	242
10.7.	<i>Practical Considerations ...</i>	245
10.8.	<i>Summary ...</i>	247
10.9.	<i>References ...</i>	248
10.10.	<i>Exercises ...</i>	250
11.	Processors	253
11.1.	<i>Introduction ...</i>	253
11.2.	<i>Tightly-Coupled Multiprocessors ...</i>	254

11.3.	<i>Priority Scheduling Disciplines ...</i>	256
11.4.	<i>Variations on Priority Scheduling ...</i>	261
11.5.	<i>FCFS Scheduling with Class-Dependent Average Service Times ...</i>	262
11.6.	<i>FCFS Scheduling with High Variability in Service Times ...</i>	263
11.7.	<i>Summary ...</i>	266
11.8.	<i>References ...</i>	268
11.9.	<i>Exercises ...</i>	270
IV.	Parameterization	273
12.	Existing Systems	274
12.1.	<i>Introduction ...</i>	274
12.2.	<i>Types and Sources of Information ...</i>	275
12.3.	<i>Customer Description ...</i>	279
12.4.	<i>Center Description ...</i>	283
12.5.	<i>Service Demands ...</i>	288
12.6.	<i>Validating the Model ...</i>	291
12.7.	<i>Summary ...</i>	293
12.8.	<i>References ...</i>	293
12.9.	<i>Exercises ...</i>	295
13.	Evolving Systems	296
13.1.	<i>Introduction ...</i>	296
13.2.	<i>Changes to the Workload ...</i>	297
13.3.	<i>Changes to the Hardware ...</i>	300
13.4.	<i>Changes to the Operating Policies and System Software ...</i>	303
13.5.	<i>Secondary Effects of Changes ...</i>	306
13.6.	<i>Case Studies ...</i>	309
13.7.	<i>Summary ...</i>	315
13.8.	<i>References ...</i>	316
13.9.	<i>Exercises ...</i>	318
14.	Proposed Systems	320
14.1.	<i>Introduction ...</i>	320
14.2.	<i>Background ...</i>	321
14.3.	<i>A General Framework ...</i>	323
14.4.	<i>Tools and Techniques ...</i>	327
14.5.	<i>Summary ...</i>	332
14.6.	<i>References ...</i>	332

V. Perspective	335
15. Extended Applications	336
15.1. <i>Introduction ...</i>	336
15.2. <i>Computer Communication Networks ...</i>	336
15.3. <i>Local Area Networks ...</i>	339
15.4. <i>Software Resources ...</i>	342
15.5. <i>Database Concurrency Control ...</i>	343
15.6. <i>Operating System Algorithms ...</i>	347
15.7. <i>Summary ...</i>	349
15.8. <i>References ...</i>	351
16. Using Queueing Network Modelling Software	354
16.1. <i>Introduction ...</i>	354
16.2. <i>Components of Queueing Network Modelling Software ...</i>	354
16.3. <i>An Example ...</i>	360
16.4. <i>Summary ...</i>	369
16.5. <i>Epilogue ...</i>	370
16.6. <i>References ...</i>	372
VI. Appendices	375
17. Constructing a Model from RMF Data	376
17.1. <i>Introduction ...</i>	376
17.2. <i>Overview of MVS ...</i>	377
17.3. <i>Overview of RMF Reports ...</i>	378
17.4. <i>Customer Description ...</i>	383
17.5. <i>Center Description ...</i>	385
17.6. <i>Service Demands ...</i>	386
17.7. <i>Performance Measures ...</i>	388
17.8. <i>Summary ...</i>	390
17.9. <i>References ...</i>	392
17.10. <i>Exercises ...</i>	393
18. An Implementation of Single Class, Exact MVA	395
18.1. <i>Introduction ...</i>	395
18.2. <i>The Program ...</i>	395
19. An Implementation of Multiple Class, Exact MVA	398
19.1. <i>Introduction ...</i>	398
19.2. <i>The Program ...</i>	398

20. Load Dependent Service Centers	403
20.1. <i>Introduction ...</i>	403
20.2. <i>Single Class Models ...</i>	405
20.3. <i>Multiple Class Models ...</i>	405
20.4. <i>Program Implementation ...</i>	407
Index	409

Part I

Preliminaries

This first part of the book provides four different sorts of background material as a prelude to our study of quantitative system performance.

In Chapter 1 we survey queueing network modelling, discussing some example applications and comparing it to more traditional approaches to computer system analysis with which the reader may be familiar.

In Chapter 2 we use case studies to explore various aspects of conducting a modelling study. Our objective is to provide some perspective on the “pieces” of the process that will be studied in the remainder of the book.

In Chapter 3 we provide a technical foundation for our work by defining a number of quantities of interest, introducing the notation that we will use in referring to these quantities, and deriving various relationships among these quantities.

In Chapter 4 we describe the inputs and the outputs of queueing network models.

Chapter 1

An Overview of Queueing Network Modelling

1.1. Introduction

Today's computer systems are more complex, more rapidly evolving, and more essential to the conduct of business than those of even a few years ago. The result is an increasing need for tools and techniques that assist in understanding the behavior of these systems. Such an understanding is necessary to provide intelligent answers to the questions of cost and performance that arise throughout the life of a system:

- *during design and implementation*
 - An aerospace company is designing and building a computer-aided design system to allow several hundred aircraft designers simultaneous access to a distributed database through graphics workstations. Early in the design phase, fundamental decisions must be made on issues such as the database accessing mechanism and the process synchronization and communication mechanism. The relative merits of various mechanisms must be evaluated prior to implementation.
 - A computer manufacturer is considering various architectures and protocols for connecting terminals to mainframes using a packet-oriented broadcast communications network. Should terminals be clustered? Should packets contain multiple characters? Should characters from multiple terminals destined for the same mainframe be multiplexed in a single packet?
- *during sizing and acquisition*
 - The manufacturer of a turn-key medical information system needs an efficient way to size systems in preparing bids. Given estimates of the arrival rates of transactions of various types, this vendor must project the response times that the system will provide when running on various hardware configurations.