

# Stochastic-Process Limits

## An Introduction to Stochastic-Process Limits And their Application to Queues

Ward Whitt

AT&T Labs - Research  
The Shannon Laboratory  
Florham Park, New Jersey

Draft  
June 13, 2001

Copyright ©info



# Preface

## 0.1. What Is This Book About?

This book is about *stochastic-process limits*, i.e., limits in which a sequence of stochastic processes converges to another stochastic process. Since the converging stochastic processes are constructed from initial stochastic processes by appropriately scaling time and space, the stochastic-process limits provide a macroscopic view of uncertainty. The stochastic-process limits are interesting and important because they generate simple approximations for complicated stochastic processes and because they help explain the statistical regularity associated with a macroscopic view of uncertainty.

This book emphasizes the continuous-mapping approach to obtain new stochastic-process limits from previously established stochastic-process limits. The continuous-mapping approach is applied to obtain stochastic-process limits for *queues*, i.e., probability models of service systems or waiting lines. These limits for queues are called *heavy-traffic limits*, because they involve a sequence of models in which the offered loads are allowed to increase towards the critical value for stability. These heavy-traffic limits generate simple approximations for complicated queueing processes under normal loading and reveal the impact of variability upon queueing performance. By focusing on the application of stochastic-process limits to queues, this book also provides an introduction to heavy-traffic stochastic-process limits for queues.

## 0.2. In More Detail

More generally, this is a book about *probability theory* – a subject which has applications to every branch of science and engineering. Probability theory can help manage a portfolio and it can help engineer a communication

network. As it should, probability theory tells us how to compute probabilities, but probability theory also has a more majestic goal: *Probability theory aims to explain the statistical regularity associated with a macroscopic view of uncertainty.*

In probability theory, there are many important ideas. But one idea might fairly lay claim to being the central idea: That idea is conveyed by the *central limit theorem*, which explains the ubiquitous bell-shaped curve: Following the giants – De Moivre, Laplace and Gauss – we have come to realize that, under regularity conditions, a sum of random variables will be approximately normally distributed if the number of terms is sufficiently large.

In the last half century, through the work of Erdős and Kac (1946, 1947), Doob (1949), Donsker (1951, 1952), Prohorov (1956), Skorohod (1956) and others, a broader view of the central limit theorem has emerged. We have discovered that there is not only statistical regularity in the  $n^{\text{th}}$  sum as  $n$  gets large, but there also is statistical regularity in the first  $n$  sums. That statistical regularity is expressed via a stochastic-process limit, i.e., a limit in which a sequence of stochastic processes converges to another stochastic process: A sequence of continuous-time stochastic processes generated from the first  $n$  sums converges in distribution to Brownian motion as  $n$  increases. That generalization of the basic central limit theorem (CLT) is known as *Donsker's theorem*. It is also called a functional central limit theorem (FCLT), because it implies convergence in distribution for many functionals of interest, such as the maximum of the first  $n$  sums. The ordinary CLT becomes a simple consequence of Donsker's FCLT, obtained by applying a projection onto one coordinate, making the ordinary CLT look like a view from Abbott's (1952) *Flatland*.

As an extension of the CLT, Donsker's FCLT is important because it has many significant applications, beyond what we would imagine knowing the CLT alone. For example, there are many applications in Statistics: Donsker's FCLT enables us to determine asymptotically-exact approximate distributions for many test statistics. The classic example is the *Kolmogorov-Smirnov statistic*, which is used to test whether data from an unknown source can be regarded as an independent sample from a candidate distribution. The stochastic-process limit identifies a relatively simple approximate distribution for the test statistic, for any continuous candidate cumulative distribution function, that can be used when the sample size is large. Indeed, early work on the Kolmogorov-Smirnov statistic by Doob (1949) and Donsker (1952) provided a major impetus for the development of the general theory of stochastic-process limits. The evolution of that story can be

seen in the books by Billingsley (1968, 1999), Csörgő and Horváth (1993), Pollard (1984), Shorack and Wellner (1986) and van der Waart and Wellner (1996).

Donsker's FCLT also has applications in other very different directions. The application that motivated this book is the application to queues: *Donsker's FCLT can be applied to establish heavy-traffic stochastic-process limits for queues.* A heavy-traffic limit for an open queueing model (with input from outside that eventually departs) is obtained by considering a sequence of queueing models, where the input load is allowed to increase toward the critical level for stability (where the input rate equals the maximum potential output rate). In such a heavy-traffic limit, the steady-state performance descriptions, such as the steady-state queue length, typically grow without bound. Nevertheless, with appropriate scaling of both time and space, there may be a nondegenerate stochastic-process limit for the entire queue-length process, which can yield useful approximations and can provide insight into system performance. The approximations can be useful even if the actual queueing systems do not experience heavy traffic. *The stochastic-process limits strip away unessential details and reveal key features determining performance.*

We are especially interested in the scaling of time and space that occurs in these heavy-traffic stochastic-process limits. It is customary to focus attention on the limit process, which serves as the approximation, but the scaling of time and space also provides important insights. For example, the scaling may reveal a *separation of time scales*, with different phenomena occurring at different time scales. In heavy-traffic limits for queues, the separation of time scales leads to unifying ideas, such as the *heavy-traffic averaging principle* (Section 2.4.2) and the *heavy-traffic snapshot principle* (Remark 5.9.1).

We obtain these many consequences of Donsker's FCLT by applying the continuous-mapping approach: Various continuous-mapping theorems imply that convergence in distribution is preserved under appropriate functions, with the simple case being a single function that is continuous. The continuous-mapping approach is much more effective with the FCLT than the CLT because many more random quantities of interest can be represented as functions of the first  $n$  partial sums than can be represented as functions of only the  $n^{\text{th}}$  partial sum. Since many heavy-traffic stochastic-process limits for queues follow from Donsker's FCLT and the continuous-mapping approach, we see that the statistical regularity revealed by the heavy-traffic limits for queues can be regarded as a consequence of the central limit theorem.

In this book we tell the story about the expanded view of the central limit theorem in more detail. We focus on stochastic-process limits, Donsker's theorem and the continuous-mapping approach. We also put life into the general theory by providing a detailed discussion of one application — queues. We give an introductory account that should be widely accessible. To help visualize the statistical regularity associated with stochastic-process limits, we perform simulations and plot stochastic-process sample paths.

However, we hasten to point out that there already is a substantial literature on stochastic-process limits, Donsker's FCLT and the continuous-mapping approach, including two editions of the masterful book by Billingsley (1968, 1999). What distinguishes the present book from previous books on this topic is our focus on *stochastic-process limits with nonstandard scaling and nonstandard limit processes*.

An important source of motivation for establishing such stochastic-process limits for queueing stochastic processes comes from evolving communication networks: Beginning with the seminal work of Leland, Taqqu, Willinger and Wilson (1994), extensive *traffic measurements* have shown that the network traffic is remarkably bursty, exhibiting complex features such as *heavy-tailed probability distributions, strong (or long-range) dependence and self-similarity*. These features present difficult engineering challenges for network design and control; e.g., see Park and Willinger (2000) and Krishnamurthy and Rexford (2001). Accordingly, a goal in our work is to gain a better understanding of these complex features and the way they affect the performance of queueing models.

To a large extent, the complex features — the heavy-tailed probability distributions, strong dependence and self-similarity — can be *defined* through their impact on stochastic-process limits. Thus, a study of stochastic-process limits, in a sufficiently broad context, is directly a study of the complex features observed in network traffic. From that perspective, it should be clear that this book is intended as a response (but not nearly a solution) to the engineering challenge posed by the traffic measurements.

We are interested in the way complex traffic affects network performance. Since a major component of network performance is congestion (queueing effects), we abstract network performance and focus on the way the complex traffic affects the performance of queues. The heavy-traffic limits show that the complex traffic can have a dramatic impact on queueing performance! We show that there are again heavy-traffic limits with these complex features, but both the scaling and the limit process may change. As in the standard case, the stochastic-process limits reveal key features determining performance.

The heavy-tailed distributions and strong dependence can lead to stochastic-process limits with *jumps in the limit process*, i.e., stochastic-process limits in which the limit process has discontinuous sample paths. The jumps have engineering significance, because they reveal sudden big changes, when viewed in a long time scale.

Much of the more technical material in the book is devoted to establishing stochastic-process limits with jumps in the limit process, but there already are books discussing stochastic-process limits with jumps in the limit process. Indeed, Jacod and Shiryaev (1987) establish many such stochastic-process limits. To be more precise, from the technical standpoint, what distinguishes this book from previous books on this topic is our focus on stochastic-process limits with *unmatched jumps* in the limit process; i.e., stochastic process limits in which the limit process has jumps unmatched in the converging processes.

For example, we may have a sequence of stochastic processes with continuous sample paths converging to a stochastic process with discontinuous sample paths. Alternatively, before scaling, we may have stochastic processes, such as queue-length stochastic processes, that move up and down by unit steps. Then, after introducing space scaling, the discontinuities are asymptotically negligible. Nevertheless, the sequence of scaled stochastic processes can converge in distribution to a limiting stochastic process with discontinuous sample paths.

Jumps are not part of Donsker's FCLT, because Brownian motion has continuous sample paths. But the classical CLT and Donsker's FCLT do not capture all possible forms of statistical regularity that can prevail. Other forms of statistical regularity emerge when the assumptions of the classical CLT no longer hold. For example, if the random variables being summed have heavy-tailed probability distributions (which here means having infinite variance), then the classical CLT for partial sums breaks down. Nevertheless, there still may be statistical regularity, but it assumes a new form. Then there is a different FCLT in which the limit process has jumps!

But the jumps in this new FCLT are *matched jumps*; each jump corresponds to an exceptionally large summand in the sums. At first glance, it is not so obvious that unmatched jumps can arise. Thus, we might regard stochastic-process limits with unmatched jumps in the limit process as pathological, and thus not worth serious attention. Part of the interest here lies in the fact that such limits, not only can occur, but routinely do occur in interesting applications. In particular, unmatched jumps in the limit process frequently occur in heavy-traffic limits for queues in the presence of heavy-tailed probability distributions. For example, in a single-server queue,

the queue-length process usually moves up and down by unit steps. Hence, when space scaling is introduced, the jumps in the scaled queue-length process are asymptotically negligible. Nevertheless, occasional exceptionally long service times can cause a rapid buildup of customers, causing the sequence of scaled queue-length processes to converge to a limit process with discontinuous sample paths. We give several examples of stochastic-process limits with unmatched jumps in the limit process in Chapter 6.

Stochastic-process limits with unmatched jumps in the limit process present technical challenges: Stochastic-process limits are customarily established by exploiting the function space  $D$  of all right-continuous  $\mathbb{R}^k$ -valued functions with left limits, endowed with the Skorohod (1956)  $J_1$  topology (notion of convergence), which is often called “the Skorohod topology.” However, that topology does not permit stochastic-process limits with unmatched jumps in the limit process.

As a consequence, to establish stochastic-process limits with unmatched jumps in the limit process, we need to use a nonstandard topology on the underlying space  $D$  of stochastic-process sample paths. Instead of the standard  $J_1$  topology on  $D$ , we use the  $M_1$  topology on  $D$ , which also was introduced by Skorohod (1956). Even though the  $M_1$  topology was introduced a long time ago, it has not received much attention. Thus, a major goal here is to provide a systematic development of the function space  $D$  with the  $M_1$  topology and associated stochastic-process limits.

It turns out the standard  $J_1$  topology is stronger (or finer) than the  $M_1$  topology, so that previous stochastic-process limits established using the  $J_1$  topology also hold with the  $M_1$  topology. Thus, while the  $J_1$  topology sometimes cannot be used, the  $M_1$  topology can almost always be used. Moreover, the extra strength of the  $J_1$  topology is rarely exploited. Thus, we would be so bold as to suggest that, *if only one topology on the function space  $D$  is to be considered, then it should be the  $M_1$  topology.*

Even though our motivation comes from queueing models and their application to describe the performance of evolving communication networks, there are many other possible applications of stochastic-process limits with jumps in the limit process. Indeed, stochastic-process limits with jumps in the limit process can arise whenever there are abrupt changes. There are natural applications to insurance, because insurance claim distributions often have heavy tails. There also are natural applications to finance, especially in the area of risk management; e.g., related to electricity derivatives. See Embrechts, Klüppelberg and Mikosch (1997), Adler, Feldman and Taqqu (1998) and Asmussen (2000).

In some cases, the fluctuations in a stochastic process are so strong



that no stochastic-process limit is possible with a limiting stochastic process having sample paths in the function space  $D$ . In order to establish stochastic-process limits involving such dramatic fluctuations, we introduce larger function spaces than  $D$ , which we call  $E$  and  $F$ . The names are chosen to suggest a natural progression starting from the space  $C$  of continuous functions and going beyond  $D$ . We define topologies on the spaces  $E$  and  $F$  analogous to the  $M_2$  and  $M_1$  topologies on  $D$ . Thus we exploit our study of the  $M$  topologies on  $D$  in this later work.

Even though the special focus here is on heavy-traffic stochastic-process limits for queues allowing unmatched jumps in the limit process, many heavy-traffic stochastic-process limits for queues have no jumps in the limit process. That is the case whenever we can directly apply the continuous-mapping approach with Donsker's FCLT. Then we deduce that reflected Brownian motion can serve as an asymptotically-exact approximation for several queueing processes in a heavy-traffic limit. In the queueing chapters we show how those classic heavy-traffic limits can be established and applied. Indeed, the book is also intended to serve as a general introduction to heavy-traffic stochastic-process limits for queues.

### 0.3. Organization of the Book

The book has fifteen chapters, which can be roughly grouped into four parts, ordered according to increasing difficulty. The level of difficulty is far from uniform: The first part is intended to be accessible with less background. It would be helpful (necessary?) to know something about probability and queues.

The *first part*, containing the first five chapters, provides an informal introduction to stochastic-process limits and their application to queues. The first part provides a broad overview, mostly without proofs, intending to complement and supplement other books, such as Billingsley (1968, 1999).

Chapter 1 uses simulation to help the reader directly experience the statistical regularity associated with stochastic-process limits. Chapter 2 discusses applications of the random walks simulated in Chapter 1. Chapter 3 introduces the mathematical framework for stochastic-process limits. Chapter 4 provides an overview of stochastic-process limits, presenting Donsker's theorem and some of its generalizations. Chapter 5 provides an introduction to heavy-traffic stochastic-process limits for queues.

The *second part*, containing Chapters 6 – 10, shows how the unmatched jumps can arise and expands the treatment of queueing models. The first chapter, Chapter 6 uses simulation to demonstrate that there should indeed

be unmatched jumps in the limit process in several examples. Chapter 7 continues the overview of stochastic-process limits begun in Chapter 4. The remaining chapters in the second part apply the stochastic-process limits, with the continuous-mapping approach, to obtain more heavy-traffic limits for queues.

The *third part*, containing Chapters 11 – 14, is devoted to the technical foundations needed to establish stochastic-process limits with unmatched jumps in the limit process. The earlier queueing chapters draw on the third part to a large extent. The queueing chapters are presented first to provide motivation for the technical foundations.

The third part begins with Chapter 11, which provides more details on the mathematical framework for stochastic-process limits, expanding upon the brief introduction in Chapter 3. Chapter 12 focuses on the function space  $D$  of right-continuous  $\mathbb{R}^k$ -valued functions with left limits, endowed with one of the nonstandard Skorohod (1956)  $M$  topologies ( $M_1$  or  $M_2$ ). As a basis for applying the continuous-mapping approach to establish new stochastic-process limits in this context, Chapter 13 shows that commonly used functions from  $D$  or  $D \times D$  to  $D$  preserve convergence with the  $M$  topologies. The third part concludes with Chapter 14, which establishes heavy-traffic limits for networks of queues.

The *fourth part*, containing Chapter 15, is more exploratory. It initiates new directions for research. Chapter 15 introduces the new spaces larger than  $D$  that can be used to express stochastic-process limits for scaled stochastic processes with even greater fluctuations.

The organization of the book is described in more detail at the end of Chapter 3, in Section 3.6.

Additional material is contained in an *Internet Supplement*. The Internet Supplement has three purposes: First, it is intended to maintain a list of corrections for errors found after the book has been published. Second, it is intended to provide supporting details, such as omitted proofs, for material in the book. Third, it is intended to provide supplementary material related to the subject of the book. Pointers to the Internet Supplement will be provided throughout the book. The initial contents of the Internet Supplement appear at the end of the book in Appendix B. The Internet Supplement is available online:

<http://www.research.att.com/~wow/supplement.html>

## 0.4. What is Missing?

Even though this book is long, it only provides introductions to stochastic-process limits and heavy-traffic stochastic-process limits for queues.

There are several different kinds of limits that can be considered for probability distributions and stochastic processes. Here we only consider central limit theorems and natural generalizations to the functions space  $D$ . We omit other kinds of limits such as large deviation principles. For large deviation principles, the continuous-mapping approach can be applied using contraction principles. Large deviations principles can be very useful for queues; see Shwartz and Weiss (1995). For a sample of other interesting probability limits (related to the Poisson clumping heuristic), see Aldous (1989).

Even though much of the book is devoted to queues, we only discuss heavy-traffic stochastic-process limits for queues. There is a large literature on queues. Nice *general introductions to queues*, at varying mathematical levels, are contained in the books by Asmussen (1987), Cooper (1982), Hall (1991), Kleinrock (1975, 1976) and Wolff (1989).

Queueing theory is intended to aid in the *performance analysis* of complex systems, such as computer, communication and manufacturing systems. We discuss performance implications of the heavy-traffic limits, but we do not discuss performance analysis in detail. Jain (1991) and Gunther (1998) discuss the performance analysis of computer systems; Bertsekas and Gallager (1987) discuss the performance analysis of communication networks; and Buzacott and Shanthikumar (1993) and Hopp and Spearman (1996) discuss the performance analysis of manufacturing systems.

Since we are motivated by evolving communication networks, we discuss queueing models that arise in that context, but we do not discuss the context itself. For background on evolving communication networks, see Keshav (1997), Kurose and Ross (2000) and Krishnamurthy and Rexford (2001). For research on communication network performance, see Park and Willinger (2000) and recent proceedings of *IEEE INFOCOM* and *ACM SIGCOMM*:

<http://www.ieee-infocom.org/2000/>

<http://www.acm.org/pubs/contents/proceedings/series/comm/>

Even within the relatively narrow domain of *heavy-traffic stochastic-process limits for queues*, we only provide an introduction. Harrison (1985) provided a previous introduction, focusing on Brownian motion and Brownian queues, the heavy-traffic limit processes rather than the heavy-traffic limits themselves. Harrison shows how martingales and the Ito stochastic calculus can be applied to calculate quantities of interest and solve control

problems. Newell (1982) provides useful perspective as well with his focus on deterministic and diffusion approximations. Harrison and Newell show that the limit processes can be used directly as approximations without considering stochastic-process limits. In contrast, we emphasize insights that can be gained from the stochastic-process limits, e.g., from the scaling.

The subject of heavy-traffic stochastic-process limits remains a very active research topic. Most of the recent interest focuses on *networks of queues with multiple classes of customers*. A principal goal is to determine good policies for scheduling and routing. That focus places heavy-traffic stochastic-process limits in the mainstream of operations research.

Multi-class queueing networks are challenging because the obvious stability criterion – having the traffic intensity be less than one at each queue – can in fact fail to be sufficient for stability; see Bramson (1994a, b). Thus, for general multi-class queueing networks, the very definition of heavy traffic is in question. For some of the recent heavy-traffic stochastic-process limits, new methods beyond the continuous-mapping approach have been required; see Bramson (1998) and Williams (1998a,b).

Discussion of the heavy-traffic approach to multi-class queueing networks, including optimization issues, can be found in the recent books by Chen and Yao (2001) and Kushner (2001), in the collections of papers edited by Yao (1994), Kelly and Williams (1995), Kelly, Zachary and Ziedins (1996), Dai (1998), McDonald and Turner (2000) and Park and Willinger (2000), and in recent papers such as Bell and Williams (2001), Harrison (2000, 2001a,b) and Kumar (2000). Hopefully, this book will help prepare readers to appreciate that important work and extend it in new directions.

# Contents

<b>Preface</b>	<b>iii</b>
0.1 What Is This Book About? . . . . .	iii
0.2 In More Detail . . . . .	iii
0.3 Organization of the Book . . . . .	ix
0.4 What is Missing? . . . . .	xi
<b>1 Experiencing Statistical Regularity</b>	<b>1</b>
1.1 A Simple Game of Chance . . . . .	1
1.1.1 Plotting Random Walks . . . . .	2
1.1.2 When the Game is Fair . . . . .	4
1.1.3 The Final Position . . . . .	10
1.1.4 Making an Interesting Game . . . . .	17
1.2 Stochastic-Process Limits . . . . .	21
1.2.1 A Probability Model . . . . .	21
1.2.2 Classical Probability Limits . . . . .	26
1.2.3 Identifying the Limit Process . . . . .	29
1.2.4 Limits for the Plots . . . . .	32
1.3 Invariance Principles . . . . .	35
1.3.1 The Range of Brownian Motion . . . . .	36
1.3.2 Relaxing the IID Conditions . . . . .	39
1.3.3 Different Step Distributions . . . . .	42
1.4 The Exception Makes the Rule . . . . .	45
1.4.1 Explaining the Irregularity . . . . .	47
1.4.2 The Centered Random Walk with $p = 3/2$ . . . . .	47
1.4.3 Back to the Uncentered Random Walk with $p = 1/2$ . . . . .	55
1.5 Summary . . . . .	60
<b>2 Random Walks in Applications</b>	<b>63</b>
2.1 Stock Prices . . . . .	63

2.2	The Kolmogorov-Smirnov Statistic . . . . .	66
2.3	A Queueing Model for a Buffer in a Switch . . . . .	70
2.3.1	Deriving the Proper Scaling . . . . .	71
2.3.2	Simulation Examples . . . . .	75
2.4	Engineering Significance . . . . .	81
2.4.1	Buffer Sizing . . . . .	81
2.4.2	Scheduling Service for Multiple Sources . . . . .	86
<b>3</b>	<b>The Framework for Stochastic-Process Limits</b>	<b>93</b>
3.1	Introduction . . . . .	93
3.2	The Space $\mathcal{P}$ . . . . .	94
3.3	The Space $D$ . . . . .	97
3.4	The Continuous-Mapping Approach . . . . .	104
3.5	Useful Functions . . . . .	106
3.6	Organization of the Book . . . . .	110
<b>4</b>	<b>A Panorama of Stochastic-Process Limits</b>	<b>115</b>
4.1	Introduction . . . . .	115
4.2	Self-Similar Processes . . . . .	116
4.2.1	General CLT's and FCLT's . . . . .	116
4.2.2	Self-Similarity . . . . .	117
4.2.3	The Noah and Joseph Effects . . . . .	120
4.3	Donsker's Theorem . . . . .	122
4.3.1	The Basic Theorems . . . . .	122
4.3.2	Multidimensional Versions . . . . .	125
4.4	Brownian Limits with Weak Dependence . . . . .	128
4.5	The Noah Effect: Heavy Tails . . . . .	132
4.5.1	Stable Laws . . . . .	133
4.5.2	Convergence to Stable Laws . . . . .	137
4.5.3	Convergence to Stable Lévy Motion . . . . .	140
4.5.4	Extreme-Value Limits . . . . .	142
4.6	The Joseph Effect: Strong Dependence . . . . .	144
4.6.1	Strong Positive Dependence . . . . .	145
4.6.2	Additional Structure . . . . .	147
4.6.3	Convergence to Fractional Brownian Motion . . . . .	150
4.7	Heavy Tails Plus Dependence . . . . .	157
4.7.1	Additional Structure . . . . .	157
4.7.2	Convergence to Stable Lévy Motion . . . . .	158
4.7.3	Linear Fractional Stable Motion . . . . .	161
4.8	Summary . . . . .	164

<b>5</b>	<b>Heavy-Traffic Limits for Fluid Queues</b>	<b>167</b>
5.1	Introduction . . . . .	167
5.2	A General Fluid-Queue Model . . . . .	169
5.2.1	Input and Available-Processing Processes . . . . .	170
5.2.2	Infinite Capacity . . . . .	171
5.2.3	Finite Capacity . . . . .	174
5.3	Unstable Queues . . . . .	177
5.3.1	Fluid Limits for Fluid Queues . . . . .	177
5.3.2	Stochastic Refinements . . . . .	181
5.4	Heavy-Traffic Limits for Stable Queues . . . . .	185
5.5	Heavy-Traffic Scaling . . . . .	191
5.5.1	The Impact of Scaling Upon Performance . . . . .	192
5.5.2	Identifying Appropriate Scaling Functions . . . . .	194
5.6	Limits as the System Size Increases . . . . .	197
5.7	Brownian Approximations . . . . .	201
5.7.1	The Brownian Limit . . . . .	201
5.7.2	The Steady-State Distribution. . . . .	203
5.7.3	The Overflow Process . . . . .	207
5.7.4	One-Sided Reflection . . . . .	210
5.7.5	First-Passage Times . . . . .	213
5.8	Planning Queueing Simulations . . . . .	215
5.8.1	The Standard Statistical Procedure . . . . .	218
5.8.2	Invoking the Brownian Approximation . . . . .	219
5.9	Heavy-Traffic Limits for Other Processes . . . . .	222
5.9.1	The Departure Process . . . . .	222
5.9.2	The Processing Time . . . . .	223
5.10	Priorities . . . . .	227
5.10.1	A Heirarchical Approach . . . . .	228
5.10.2	Processing Times . . . . .	230
<b>6</b>	<b>Unmatched Jumps in the Limit Process</b>	<b>233</b>
6.1	Introduction . . . . .	233
6.2	Linearly Interpolated Random Walks . . . . .	235
6.2.1	Asymptotic Equivalence with $M_1$ . . . . .	236
6.2.2	Simulation Examples . . . . .	237
6.3	Heavy-Tailed Renewal Processes . . . . .	241
6.3.1	Inverse Processes . . . . .	241
6.3.2	The Special Case with $m = 1$ . . . . .	244
6.4	A Queue with Heavy-Tailed Distributions . . . . .	250
6.4.1	The Standard Single-Server Queue . . . . .	251

6.4.2	Heavy-Traffic Limits . . . . .	252
6.4.3	Simulation Examples . . . . .	254
6.5	Rare Long Service Interruptions . . . . .	263
6.6	Time-Dependent Arrival Rates . . . . .	267
<b>7</b>	<b>More Stochastic-Process Limits</b>	<b>273</b>
7.1	Introduction . . . . .	273
7.2	Central Limit Theorem for Processes . . . . .	274
7.2.1	Hahn's Theorem . . . . .	274
7.2.2	A Second Limit . . . . .	279
7.3	Counting Processes . . . . .	283
7.3.1	CLT Equivalence . . . . .	284
7.3.2	FCLT Equivalence . . . . .	285
7.4	Renewal-Reward Processes . . . . .	289
<b>8</b>	<b>Fluid Queues with On-Off Sources</b>	<b>295</b>
8.1	Introduction . . . . .	295
8.2	A Fluid Queue Fed by On-Off Sources . . . . .	298
8.2.1	The On-Off Source Model . . . . .	298
8.2.2	Simulation Examples . . . . .	301
8.3	Heavy-Traffic Limits for the On-Off Sources . . . . .	307
8.3.1	A Single Source . . . . .	307
8.3.2	Multiple Sources . . . . .	310
8.3.3	$M/G/\infty$ Sources . . . . .	314
8.4	Brownian Approximations . . . . .	316
8.4.1	The Brownian Limit . . . . .	316
8.4.2	Model Simplification . . . . .	319
8.5	Stable-Lévy Approximations . . . . .	321
8.5.1	The RSLM Heavy-Traffic Limit . . . . .	321
8.5.2	The Steady-State Distribution . . . . .	325
8.5.3	Numerical Comparisons . . . . .	328
8.6	Second Stochastic-Process Limits . . . . .	330
8.6.1	$M/G/1/K$ Approximations . . . . .	331
8.6.2	Limits for Limit Processes . . . . .	337
8.7	Reflected Fractional Brownian Motion . . . . .	339
8.7.1	An Increasing Number of Sources . . . . .	339
8.7.2	Gaussian Input . . . . .	340
8.8	Reflected Gaussian Processes . . . . .	343



<b>9</b>	<b>Single-Server Queues</b>	<b>347</b>
9.1	Introduction . . . . .	347
9.2	The Standard Single-Server Queue . . . . .	349
9.3	Heavy-Traffic Limits . . . . .	353
9.3.1	The Scaled Processes . . . . .	353
9.3.2	Discrete-Time Processes . . . . .	356
9.3.3	Continuous-Time Processes . . . . .	359
9.4	Superposition Arrival Processes . . . . .	364
9.5	Split Processes . . . . .	368
9.6	Brownian Approximations . . . . .	370
9.6.1	Variability Parameters . . . . .	371
9.6.2	Models with More Structure . . . . .	374
9.7	Very Heavy Tails . . . . .	378
9.7.1	Heavy-Traffic Limits . . . . .	379
9.7.2	First Passage to High Levels . . . . .	380
9.8	An Increasing Number of Arrival Processes . . . . .	383
9.8.1	Iterated and Double Limits . . . . .	383
9.8.2	Separation of Time Scales . . . . .	389
9.9	Approximations for Queueing Networks . . . . .	393
9.9.1	Parametric-Decomposition Approximations . . . . .	393
9.9.2	Approximately Characterizing Arrival Processes . . . . .	398
9.9.3	A Network Calculus . . . . .	399
9.9.4	Exogenous Arrival Processes . . . . .	406
9.9.5	Concluding Remarks . . . . .	407
<b>10</b>	<b>Multi-Server Queues</b>	<b>411</b>
10.1	Introduction . . . . .	411
10.2	Queues with Multiple Servers . . . . .	412
10.2.1	A Queue with Autonomous Service . . . . .	412
10.2.2	The Standard $m$ -Server Model . . . . .	415
10.3	Infinitely Many Servers . . . . .	419
10.3.1	Heavy-Traffic Limits . . . . .	420
10.3.2	Gaussian Approximations . . . . .	424
10.4	An Increasing Number of Servers . . . . .	428
10.4.1	Infinite-Server Approximations . . . . .	428
10.4.2	Heavy-Traffic Limits for Delay Models . . . . .	430
10.4.3	Heavy-Traffic Limits for Loss Models . . . . .	433
10.4.4	Planning Simulations of Loss Models . . . . .	435

<b>11 More on the Mathematical Framework</b>	<b>441</b>
11.1 Introduction . . . . .	441
11.2 Topologies . . . . .	442
11.2.1 Definitions . . . . .	442
11.2.2 Separability and Completeness . . . . .	445
11.3 The Space $\mathcal{P}$ . . . . .	447
11.3.1 Probability Spaces . . . . .	447
11.3.2 Characterizing Weak Convergence . . . . .	448
11.3.3 Random Elements . . . . .	450
11.4 Product Spaces . . . . .	453
11.5 The Space $D$ . . . . .	457
11.5.1 $J_2$ and $M_2$ Metrics . . . . .	457
11.5.2 The Four Skorohod Topologies . . . . .	460
11.5.3 Measurability Issues . . . . .	462
11.6 The Compactness Approach . . . . .	464
<b>12 The Space <math>D</math></b>	<b>471</b>
12.1 Introduction . . . . .	471
12.2 Regularity Properties of $D$ . . . . .	472
12.3 Strong and Weak $M_1$ Topologies . . . . .	474
12.3.1 Definitions . . . . .	475
12.3.2 Metric Properties . . . . .	477
12.3.3 Properties of Parametric Representations . . . . .	480
12.4 Local Uniform Convergence at Continuity Points . . . . .	483
12.5 Alternative Characterizations of $M_1$ Convergence . . . . .	486
12.5.1 $SM_1$ Convergence . . . . .	486
12.5.2 $WM_1$ Convergence . . . . .	491
12.6 Strengthening the Mode of Convergence . . . . .	493
12.7 Characterizing Convergence with Mappings . . . . .	494
12.8 Topological Completeness . . . . .	497
12.9 Non-Compact Domains . . . . .	498
12.10 Strong and Weak $M_2$ Topologies . . . . .	501
12.11 Alternative Characterizations of $M_2$ Convergence . . . . .	504
12.11.1 $M_2$ Parametric Representations . . . . .	504
12.11.2 $SM_2$ Convergence . . . . .	505
12.11.3 $WM_2$ Convergence . . . . .	507
12.11.4 Additional Properties of $M_2$ Convergence . . . . .	509
12.12 Compactness . . . . .	511

<b>13 Useful Functions</b>	<b>515</b>
13.1 Introduction . . . . .	515
13.2 Composition . . . . .	516
13.3 Composition with Centering . . . . .	520
13.4 Supremum . . . . .	525
13.5 One-Dimensional Reflection . . . . .	529
13.6 Inverse . . . . .	532
13.6.1 The Standard Topologies . . . . .	533
13.6.2 The $M'_1$ Topology . . . . .	536
13.6.3 First Passage Times . . . . .	538
13.7 Inverse with Centering . . . . .	540
13.8 Counting Functions . . . . .	547
<b>14 Queueing Networks</b>	<b>551</b>
14.1 Introduction . . . . .	551
14.2 The Multidimensional Reflection Map . . . . .	555
14.2.1 A Special Case . . . . .	555
14.2.2 Definition and Characterization . . . . .	556
14.2.3 Continuity and Lipschitz Properties . . . . .	561
14.3 The Instantaneous Reflection Map . . . . .	570
14.3.1 Definition and Characterization . . . . .	571
14.3.2 Implications for the Reflection Map . . . . .	578
14.4 Reflections of Parametric Representations . . . . .	581
14.5 $M_1$ Continuity Results and Counterexamples . . . . .	584
14.5.1 $M_1$ Continuity Results . . . . .	584
14.5.2 Counterexamples . . . . .	587
14.6 Limits for Stochastic Fluid Networks . . . . .	590
14.6.1 Model Continuity . . . . .	592
14.6.2 Heavy-Traffic Limits . . . . .	593
14.7 Queueing Networks with Service Interruptions . . . . .	596
14.7.1 Model Definition . . . . .	596
14.7.2 Heavy-Traffic Limits . . . . .	600
14.8 The Two-Sided Regulator . . . . .	607
14.8.1 Definition and Basic Properties . . . . .	608
14.8.2 With the $M_1$ Topologies . . . . .	612
14.9 Chapter Notes . . . . .	615

<b>15 The Spaces <math>E</math> and <math>F</math></b>	<b>619</b>
15.1 Introduction . . . . .	619
15.2 Three Time Scales . . . . .	620
15.3 More Complicated Oscillations . . . . .	624
15.4 The Space $E$ . . . . .	629
15.5 Characterizations of $M_2$ Convergence in $E$ . . . . .	634
15.6 Convergence to Extremal Processes . . . . .	637
15.7 The Space $F$ . . . . .	641
15.8 Queueing Applications . . . . .	643
<b>16 Bibliography</b>	<b>649</b>
<b>A Regular Variation</b>	<b>693</b>
<b>B Contents of the Internet Supplement</b>	<b>697</b>