Giuseppe Serazzi

# Performance Engineering

## Learning Through Applications Using JMT

Springer

Performance Engineering

Giuseppe Serazzi

# Performance Engineering

Learning Through Applications Using JMT

Giuseppe Serazzi
Milano, Pavia, Italy

*In memory of*
*Larry Dowdy, Martin Reiser, Paul Schweitzer, Kennet Sevcik*
*before the night falls*

# Preface

This open-access book aims to improve users' skills needed to implement models for performance evaluation of digital infrastructures. Models are widely used in any branch of engineering. Unfortunately, their use for performance evaluation of computing infrastructures is pretty much limited to modeling specialists and not to their end-users, who have complete knowledge of the analyzed phenomena. Among the reasons for this limitation, there is the intrinsic complexity of the modeling process, which cannot be fully learned with the academic approach alone, and the frequent use of unnecessary mathematical details, which typically create a fog shield that hides the key features of the models. Furthermore, it is known that to increase the ability to build reliable models, it is necessary to accumulate experiences that can only be learned through trial-and-error work by solving problems of different difficulties.

Based on these considerations, we tried to keep this book as simple as possible by adopting the following guidelines. On one side, we present a collection of modeling studies of increasing complexity, describing the assumptions made and motivating the decisions taken. Readers are introduced to the modeling process gradually, learning the basic concepts step-by-step as they go through the case studies. On the other side, we try to avoid superfluous exposure to mathematical concepts. For interested readers, we reported in *Appendix* some basic notions that may be useful to know.

Among the various techniques used in performance evaluation modeling, we will use *Queueing Networks*, possibly integrated with *Petri Nets* when the characteristics of the models require it. This type of model provides a good balance between the accuracy of results, complexity, and parameterization effort, for a large variety of problems. Analytical, simulation, and asymptotic techniques are applied to solve the models.

The book is structured in six Chapters and an Appendix. Chapter 1 is focused on the description of the model building process. The input parameters, the output metrics, and the operational laws are illustrated. The most important steps to building models to be solved with simulation and analytical techniques are reviewed. In the following Chapters, fifteen case studies of increasing complexity covering different aspects of performance evaluation are described. We believe that readers could benefit

from analyzing these models by focusing on the abstraction process applied to their design. In Chaps. 2 and 3, models of systems with homogeneous and heterogeneous workloads are presented, respectively. The problem of bottleneck identification and performance optimization is addressed for both types of workloads. Chapter 4 is devoted to the analysis of the impact of variability of the traffic of requests and service demands on throughput and response time. Chapter 5 focuses on parallel computing and describes the influence of different synchronization policies on performance. Chapter 6 presents four case studies derived from real-life scenarios: a surveillance system, an architecture that autoscales for load fluctuations, a web app workflow simulation, and a crowd computing platform. The autoscaler model consists of Queueing Networks and Petri Nets integrated, i.e., it is a *multi-formalism* model.

The `Java Modelling Tools (JMT)`, a *open source* suite of six tools for performance engineering and capacity planning using Queueing Networks and Petri Nets, were applied to build and solve the models. Details on `JMT`, which can be downloaded from `http://jmt.sourceforge.net`, can be found in [8]. `JMT` is a project coordinated and co-developed by Politecnico di Milano (G. Serazzi) and Imperial College London (G. Casale).

This book is intended as a text for courses in performance evaluation and modeling for graduate and senior-level computer science students. Researchers and practitioners whose work is related to performance evaluation of computer infrastructures will find it useful as a reference text. It can be used also as a supporting text for courses in disciplines outside of computer science that require the use of modeling to evaluate the performance of their applications.

We hope you will find the *learning through applications* approach followed in this book useful for your work, and apologize in advance for the mistakes you will find. The author cannot be considered responsible for errors that you may introduce in your work due to the content of this book.

Milano, Pavia, Italy                                                            Giuseppe Serazzi
May 2023

# Acknowledgments

# Contents