

Systems & Control: Foundations & Applications

R. Srikant

The Mathematics of Internet Congestion Control

Birkhäuser

R. Srikant

The Mathematics of Internet Congestion Control



Birkhäuser
Boston • Basel • Berlin

R. Srikant
Coordinated Science Laboratory and
Department of General Engineering
University of Illinois
Urbana, IL 61801
USA

Library of Congress Cataloging-in-Publication Data

Srikant, Rayadurgam.

The mathematics of Internet congestion control / Rayadurgam Srikant.

p. cm. – (Systems and control : foundations and applications)

Includes bibliographical references and index.

1. Internet—Mathematical models. 2.

Telecommunication—Traffic—Management—Mathematics. I. Title. II. Systems & control.

TK5105.875.I57S685 2003

004.67'8'015118—dc22

2003065233

CIP

AMS Subject Classifications: 90C25, 93D2, 93A15, 93D30, 93E15, 90B18, 90B10, 90B15

ISBN 0-8176-3227-1

Printed on acid-free paper.

©2004 Birkhäuser Boston

Birkhäuser 

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Birkhäuser Boston, c/o Springer-Verlag New York, Inc., 175 Fifth Avenue, New York, NY 10010, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed in the United States of America. (SB)

9 8 7 6 5 4 3 2 1

SPIN 10925164

Birkhäuser is part of *Springer Science+Business Media*

www.birkhauser.com



Systems & Control: Foundations & Applications

Series Editor

Tamer Başar, University of Illinois at Urbana-Champaign

Editorial Board

Karl Johan Åström, Lund Institute of Technology, Lund, Sweden

Han-Fu Chen, Academia Sinica, Beijing

William Helton, University of California, San Diego

Alberto Isidori, University of Rome (Italy) and

Washington University, St. Louis

Petar V. Kokotović, University of California, Santa Barbara

Alexander Kurzhanski, Russian Academy of Sciences, Moscow
and University of California, Berkeley

H. Vincent Poor, Princeton University

Mete Soner, Koç University, Istanbul

*To
Amma, Appa
Susie, Katie, and Jenny*

Preface

The Transmission Control Protocol (TCP) was introduced in the 1970s to facilitate reliable file transfer in the Internet. However, there was very little in the original TCP to control congestion in the network. If several users started transferring files over a single bottleneck link at a total rate exceeding the capacity of the link, then packets had to be dropped. This resulted in retransmissions of lost packets, which again led to more lost packets. This phenomenon known as *congestion collapse* occurred many times in the mid-1980s, prompting the development of a congestion control mechanism for TCP by Van Jacobson. By all accounts, this has been a remarkably successful algorithm steering the Internet through an era of unprecedented global expansion. However, with access speeds to the Internet having grown by several orders of magnitude over the past decade and round-trip times increasing due to the global nature of the Internet, there is a need to develop a more scalable mechanism for Internet congestion control. Loosely speaking, by scalability, we mean that the protocol should exhibit a provably good behavior which is unaffected by the number of nodes in the Internet, the capacities of the links, and the RTTs (round-trip times) involved. The purpose of this book is to provide an introduction to the significant progress in the mathematical modelling of congestion control which was initiated by work of Frank Kelly in the mid-1990s, and further developed by many researchers since then.

The book draws upon results from three widely-used topics to model the Internet: convex optimization, control theory and probability. It would be difficult to read this book without a level of knowledge equivalent to that of a first undergraduate course in each of these three topics. On the other hand, it does not require much more than a first-level course in these topics to be able to understand most of the material presented in the book. At the end of most chapters, a brief appendix is included with the intent of providing some background material that would be useful in understanding the development in the main body of the book.

While the mathematical modelling of Internet congestion control is a fairly recent topic, due to the importance of the topic, it has attracted a large

number of researchers who have made important contributions to this subject. From among this vast body of literature, I have chosen to focus in this book on those approaches which are rooted in the view that congestion control is a mechanism for resource allocation in a network. Even with this framework, there is a large body of work and I have primarily chosen to emphasize scalable, decentralized mechanisms in this book.

Many people have directly or indirectly contributed to my understanding of the topic of mathematical modelling of the Internet. I wish to thank Frank Kelly whose seminal work was the impetus to all the topics discussed in the book. I have learned a lot about congestion control both from reading his papers as well as from interactions with him at conferences and through email discussions. It is a pleasure to acknowledge and thank Tamer Başar for earlier collaboration on the ATM ABR service and more recent collaboration on the primal-dual algorithm for TCP congestion control. I would also like to acknowledge the work of many present and past graduate students who have contributed to many of the results presented in this book. They include Srisankar Kunniyur, Sanjay Shakkottai, Supratim Deb, Damien Polis, Srinivas Shakkottai, Ashvin Lakshmikantha, Julian Shao Liu and Lei Ying. I also gratefully acknowledge the help of Supratim in generating some of the figures used in this book and providing comments on Chapter 8, and Julian for proof-reading the entire manuscript. I would also like to thank Eitan Altman, Carolyn Beck, Geir Dullerud, Peter Key, A.J. Ganesh, Don Towsley and Chris Holot for collaborations that have contributed to my understanding of the topics in the book. Finally, I thank Bruce Hajek for always being available to discuss any mathematical problem on any topic. My discussions with him have greatly shaped the direction of much of my research in general.

Urbana, IL
October 2003

R. Srikant

Contents

Preface	vii
1 Introduction	1
2 Resource Allocation	7
2.1 Resource allocation as an optimization problem	8
2.2 A general class of utility functions	13
2.3 Appendix: Convex optimization	17
3 Congestion Control: A decentralized solution	23
3.1 Primal algorithm	23
3.2 Dual algorithm	27
3.3 Exact penalty functions	29
3.4 Primal-dual approach	32
3.5 Other variations in the primal approach	33
3.5.1 Exponentially averaged rate feedback	33
3.5.2 One-bit marking in the primal method	34
3.6 REM: A one-bit marking scheme	35
3.7 Multipath routing	36
3.8 Multirate multicast congestion control	37
3.9 A pricing interpretation of proportional fairness	43
3.10 Appendix: Lyapunov stability	46
4 Relationship to Current Internet Protocols	49
4.1 Window flow control	49
4.2 Jacobson's adaptive window flow control algorithm	50
4.2.1 TCP and resource allocation	59
4.3 TCP-Vegas	61
4.4 Random Early Detection (RED)	63
4.5 Explicit Congestion Notification (ECN)	65
4.6 High-throughput TCP	65

5	Linear Analysis with Delay: The single link case	67
5.1	Single TCP-Reno source with droptail	67
5.2	Multiple TCP sources with identical RTTs	69
5.3	TCP-Reno and RED	72
5.4	Proportionally-fair controller	75
5.5	High-throughput TCP	76
5.6	Dual algorithm	77
5.7	Primal-dual algorithm	79
5.8	Appendix: The Nyquist criterion	80
6	Linear Analysis with Delay: The network case	83
6.1	Primal controllers	83
6.1.1	Proportionally-fair controller	83
6.1.2	High-throughput TCP with rate-based feedback	87
6.1.3	Exponentially smoothed rate feedback	90
6.1.4	High-throughput TCP with probabilistic marking	92
6.1.5	A general nonlinear increase/decrease algorithm	94
6.1.6	High-throughput TCP and AVQ	96
6.2	Dual algorithm	101
6.2.1	Fair dual algorithm	103
6.3	Primal-dual algorithm	104
6.4	Appendix: Multivariable Nyquist criterion	108
7	Global Stability for a Single Link and a Single Flow	109
7.1	Proportionally-fair controller over a single link	109
8	Stochastic Models and their Deterministic Limits	117
8.1	Deterministic limit for proportionally-fair controllers	118
8.2	Individual source dynamics	125
8.3	Price feedback	128
8.4	Queue-length-based marking	129
8.5	TCP-type congestion controllers	131
8.6	Appendix: The weak law of large numbers	133
9	Connection-level Models	135
9.1	Stability of weighted proportionally-fair controllers	135
9.2	Priority resource allocation	139
10	Real-time Sources and Distributed Admission Control	141
10.1	Resource sharing between elastic and inelastic users	142
10.2	Probing and distributed admission control	144
10.3	A simple model for queueing at the link buffer	146
10.4	Appendix: Diffusion approximation	148
10.4.1	Brownian motion through a queue	149
10.4.2	A lower bound for $\lim_{t \rightarrow \infty} P(q(t) > x)$	151

10.4.3 Computing the steady-state density using PDE's	152
11 Conclusions	155
References	157
Index	163

Introduction

The Internet has evolved from a loose federation of networks used primarily in academic institutions, to a global entity which has revolutionized communication, commerce and computing. Early in the evolution, it was recognized that unrestricted access to the Internet resulted in poor performance in the form of low network utilization and high packet loss rates. This phenomenon known as congestion collapse, led to the development of the first congestion control algorithm for the Internet [39]. The basic idea behind the algorithm was to detect congestion in the network through packet losses. Upon detecting a packet loss, the source reduces its transmission rate; otherwise, it increases the transmission rate. The original algorithm has undergone many minor, but important changes, but the essential features of the algorithm used for the increase and decrease phases of the algorithm have not changed through the various versions of TCP, such as TCP-Tahoe, Reno, NewReno, SACK [17, 54]. An exception to this is the TCP Vegas algorithm which uses queueing delay in the network as the indicator of congestion, instead of packet loss [14]. One of the goals of the book is to understand the dynamics of Jacobson's algorithm through simple mathematical models, and to develop tools and techniques that will improve the algorithm and make it scalable for networks with very large capacities, very large numbers of users, and very large round-trip times.

In parallel with the development of congestion control algorithms for the Internet, congestion control was studied for other data networks of the time. A simple, yet popular, mathematical model for allocating resources in a fair manner between two users sharing a single link was developed by Chiu and Jain in [16]. This was an early algorithm that recognized the connection between congestion control and resource allocation. We present the model and its analysis below.

Consider two sources accessing a link that can serve packets at the rate c packets per sec. Let x_i be the rate at which source i is injecting packets into the network. Suppose that the link provides feedback to the source indicating whether the total arrival rate at the link ($x_1 + x_2$) is greater than the link capacity or not. Thus, the feedback is $I_{x_1+x_2 > c}$, the indicator function of

the event $x_1 + x_2 > c$. The sources respond to this congestion indication by adjusting their rates as follows: for each i , x_i evolves according to the differential equation

$$\dot{x}_i = I_{x_1+x_2 \leq c} - \beta x_i I_{x_1+x_2 > c}. \quad (1.1)$$

To see how this system behaves, define $y = x_1 - x_2$. Then, $y(t)$ evolves according to the differential equation

$$\dot{y} = -\beta y I_{x_1+x_2 > c}.$$

When $x_1 + x_2 \leq c$, then clearly y does not change with time, thus, $x_1 - x_2$ remains a constant. However, from (1.1), x_1 and x_2 increase. To consider the behavior of x_1 and x_2 when $x_1 + x_2 \geq c$, let $V(y) = y^2$. Then,

$$\dot{V} = -2y\dot{y} = -2y^2 I_{x_1+x_2 > c}.$$

Thus, when $x_1 + x_2 > c$, $\dot{V} < 0$ unless $y = 0$. When $x_1 + x_2 \leq c$, it can be seen from (1.1) that both x_1 and x_2 increase, while $x_1 - x_2$ remains a constant. Thus, with a little bit more work, one can conclude that, as $t \rightarrow \infty$, $x_1 + x_2 \rightarrow c$, and $y \rightarrow 0$. In other words, in steady-state, the link is shared equally between the two sources and the link is fully utilized. The Chiu-Jain algorithm identifies several features of congestion management which are of interest to us:

- Congestion control: The sources control their rates x_1 and x_2 depending on the level of congestion in the network. If the arrival rate at the link is too large, then the sources decrease their transmission rates and if it is too small, then the sources increase their rates.
- Congestion feedback: The network (in this case, a single link) participates in the congestion management process by providing feedback in the form of $I_{x_1+x_2 > c}$. Note that the amount of feedback required is minimal; it requires only one bit of information from the link: whether the arrival rate exceeds the capacity at the link or not. The link does not even have to actively participate in the feedback process. If we assume that packets are dropped when the arrival rate exceeds capacity, then the receivers can detect lost packets and inform the source that there is congestion in the network.
- Fairness in resource allocation: The goal of the congestion control algorithm can be viewed as driving the system towards a fair operating point, which in this case corresponds to each user getting half of the available bandwidth.
- Utilization: The link is fully utilized, i.e., at equilibrium, the arrival rate is equal to the available capacity.
- Decentralization: The congestion controllers are decentralized. Each controller needs only one bit of information from the network, but requires no communication with the other controller(s).

The discrete-time version of (1.1) is given by

$$x_i(k+1) = x_i(k) + \delta I_{x_1(k)+x_2(k) \leq c} - \beta \delta I_{x_1(k)+x_2(k) > c},$$

for some small $\delta > 0$. For $c = 1$, $\delta = 0.05$ and $\beta = 1/2$, a plot of the evolution of the user rates is shown in Figure 1.1 starting from the initial condition $x_1 = 0.3$ and $x_2 = 0.1$.

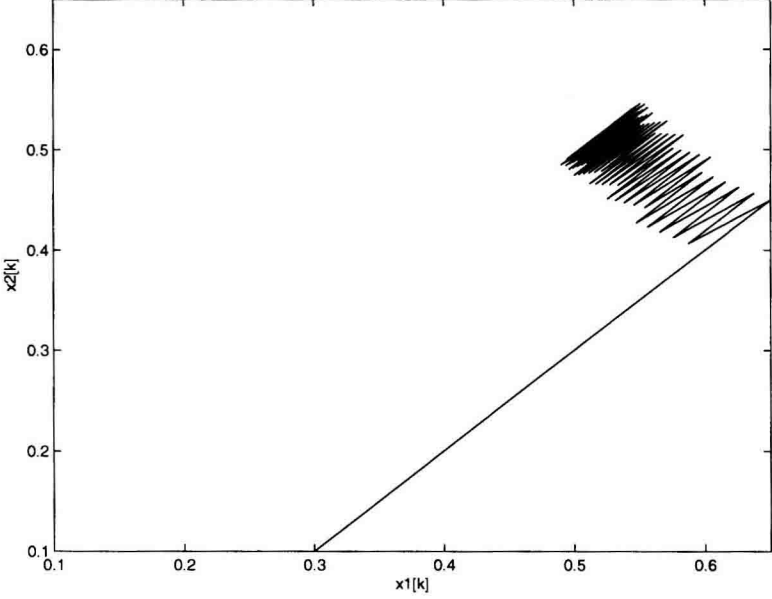


Fig. 1.1. Rate evolution using the Chiu-Jain algorithm. Starting from the point $(0.3, 0.1)$, the system moves towards the point $(0.5, 0.5)$

Next, let us consider the system shown in Figure 1.2. In this system, again there are two sources accessing a single link. However, a packet from source i takes $d_1(i)$ time units to reach the link and it takes $d_2(i)$ time units for the feedback from the link to reach source i . Thus, the dynamics of the source i are given by

$$\dot{x}_i = I_{x(t-d_2(i)) \leq c} - \beta x_i I_{x(t-d_2(i)) > c},$$

where $x(t)$ is the total arrival rate at the link at time t and is given by

$$x(t) = x_1(t - d_1(1)) + x_2(t - d_1(2)).$$

Following [13], let us consider the evolution of source i 's rate at some time $t + d_2(i)$:

$$\begin{aligned} \dot{x}_i(t + d_2(i)) &= I_{x_1(t-d_1(1))+x_2(t-d_1(2)) \leq c} \\ &\quad - \beta x_i(t + d_2(i)) I_{x_1(t-d_1(1))+x_2(t-d_1(2)) > c}, \end{aligned}$$

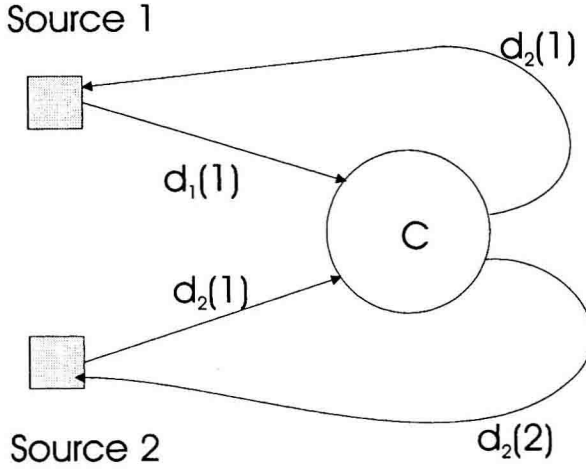


Fig. 1.2. A link with sources having delays in the forward and reverse paths

and define

$$y = x_1(t + d_2(1)) - x_2(t + d_2(2)).$$

If $V(y) = y^2$, then

$$\dot{y} = -2\beta y I_{x_1(t-d_1(1))+x_2(t-d_1(2))>c},$$

from which we can conclude as in the delay-free case that $y \rightarrow 0$ as $t \rightarrow \infty$. In the presence of delay, note that even in steady-state, the rate allocation is not fair at every time instant, but is fair if one compares time-shifted versions of the two rates. In other words, on average, each user gets its fair share of the bandwidth. A discrete version of the algorithm with time step $\delta = 0.005$, $\beta = 0.5$, $d_1(1) = d_2(1) = 2$ discrete time steps, $d_1(2) = d_2(2) = 5$ time steps was simulated and the results are shown in Figure 1.3.

It is difficult to generalize this algorithm and to establish its convergence for more general topology networks. Further, it is not immediately clear if this simple notion of fairness can be generalized to networks with more than one link. However, this algorithm was historically influential in providing the motivation for the congestion control mechanism in [39] which will be studied in detail in a later chapter. The one-bit feedback mechanism proposed by Chiu and Jain and the work reported in [87] are also precursors to the later one-bit feedback schemes such as RED [27].

After the introduction of the congestion control algorithm for TCP, several researchers developed simple mathematical models that led to a better understanding of the dynamics of this algorithm (see, for example, [65, 73, 56]). Building on a resource allocation model by Kelly [48], the seminal paper by Kelly, Mauloo and Tan [52] presented the first mathematical model and analysis of the behavior of congestion control algorithms for general topology networks. Since then, there have been significant developments in the modelling

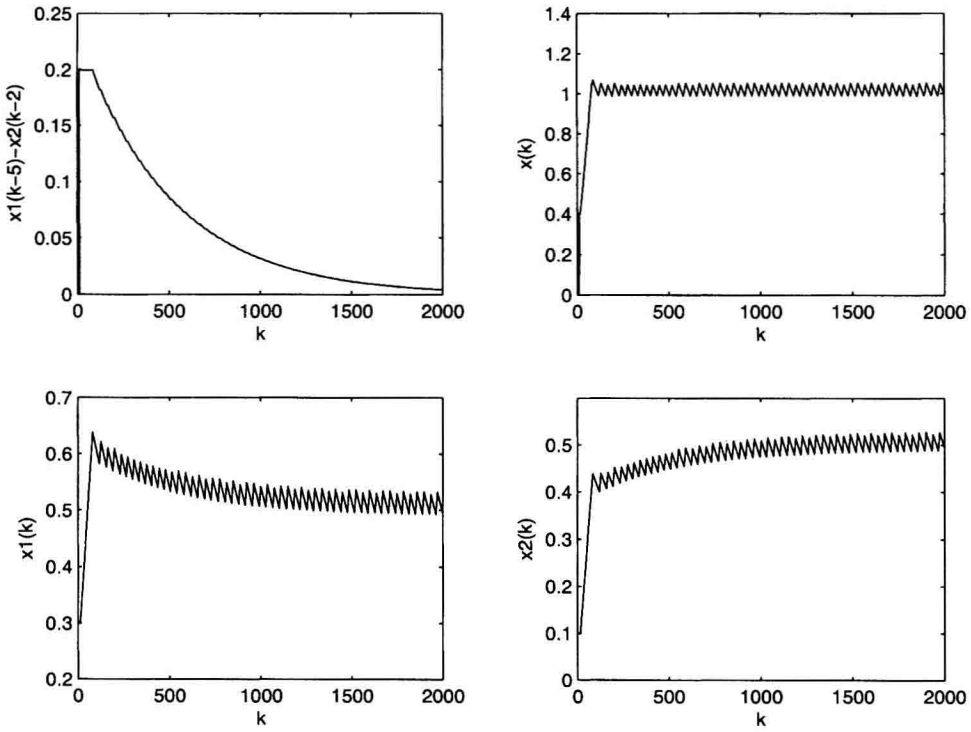


Fig. 1.3. Simulation results of the Chiu-Jain algorithm with feedback delays

of Internet congestion control using tools from operations research and control theory. The goal of this book is to provide a comprehensive introduction to this remarkable progress in the development of an Internet congestion control theory.

