

Fairness in Computer Networks

Inauguraldissertation
zur Erlangung des akademischen Grades
eines Doktors der Naturwissenschaften
der Universität Mannheim

vorgelegt von
Diplom Wirtschaftsinformatiker Robert R. Denda
aus Koblenz

Mannheim, 2003

Dekan: Professor Dr. Jürgen Potthoff, Universität Mannheim
Referent: Professor Dr. Wolfgang Effelsberg, Universität Mannheim
Korreferent: Professor Dr. Stefan Fischer, Technische Universität Braunschweig

Tag der mündlichen Prüfung: 15. Januar 2004

Berichte aus der Telematik

Robert Denda

Fairness in Computer Networks

D 180 (Diss. Universität Mannheim)

Shaker Verlag
Aachen 2004

Bibliographic information published by Die Deutsche Bibliothek

Die Deutsche Bibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data is available in the internet at <http://dnb.ddb.de>.

Zugl.: Mannheim, Univ., Diss., 2004

Copyright Shaker Verlag 2004

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publishers.

Printed in Germany.

ISBN 3-8322-2581-1

ISSN 0948-700X

Shaker Verlag GmbH • P.O. BOX 101818 • D-52018 Aachen

Phone: 0049/2407/9596-0 • Telefax: 0049/2407/9596-9

Internet: www.shaker.de • eMail: info@shaker.de

Kurzfassung

Dienstgüte (Quality-of-Service) war im Rechnernetzwerkbereich eines der erfolgreichsten Forschungsgebiete der letzten zehn Jahre. Basierend auf den gewonnenen Forschungsergebnissen bedarf es nun hauptsächlich eines kontinuierlichen Standardisierungsprozesses und eines weitergehenden kommerziellen Einsatzes, um der sich entwickelnden Nachfrage nach Dienstgüte moderner Netzwerkanwendungen nachzukommen.

Das zur Dienstgüte komplementäre Gebiet *Fairness* wurde bisher jedoch nicht ausreichend tief betrachtet, wozu die vorliegende Dissertation einen wichtigen Forschungsbeitrag liefert.

Wir stellen die Verbindung zwischen den theoretischen Ergebnissen aus den zugehörigen Bereichen der Volkswirtschaftstheorie und den entsprechenden Anforderungen von Rechnernetzen her. Fairnessdefinitionen werden aufgezeigt und ihr Nutzen für den Netzwerkbereich evaluiert. Des Weiteren beinhaltet die vorliegende Dissertation eine erste komplett Klassifizierung bestehender Arbeiten über Fairness in Rechnernetzen. Dafür haben wir Fairnessaspekte in die beiden Bereiche *Makro-Fairness* und *Mikro-Fairness* aufgeteilt.

Makro-Fairness bezieht sich auf eine allgemein faire Verteilung der Netzwerkressourcen, von der alle Nutzer und Anwendungen profitieren, wohingegen *Mikro-Fairness* sich auf Fairness im engeren Sinne für spezifische Anwendungen in Szenarien mit ausgeprägtem Wettbewerb, wie z.B. Echtzeithandelssysteme und Online-Auktionen, bezieht.

Im Bereich der Makro-Fairness liefert die vorliegende Dissertation zwei wichtige Forschungsbeiträge: einen neuartigen Architekturvorschlag, *iSSD*, der nutzerbasierte Fairness mit ähnlich niedriger Komplexität wie DiffServ erreicht, und ein neues TCP-freundliches Congestion Control Protokoll, *MPLBcc*. Die Fairnesseigenschaften der iSSD Architektur werden im Rahmen dieser Dissertation sowohl durch eine formale Analyse als auch durch Netzwerksimulationen bestätigt. Der zweite Forschungsbeitrag für den Bereich Makro-Fairness, MPLBcc, liefert für den speziellen Fall mehrerer gleichzeitiger Übertragungspfade bessere Ergebnisse als bestehende TCP-freundliche Algorithmen. Sowohl Simulationsergebnisse als auch Ergebnisse des Feldeinsatzes einer implementierten Variante von MPLBcc, die im Rahmen eines größeren Projektes für die Übertragung von digitalem Video über den Kabelnetzrückkanal eingesetzt wurde, werden aufgezeigt.

Bezüglich Mikro-Fairness, werden von uns entwickelte und getestete neue Algorithmen vorgestellt, die eine faire Übertragung von Netzwerkströmen des Typs *group-delay constrained traffic* ermöglichen, was besonders für wettbewerbsorientierte Anwendungen, wie z.B. Aktienhandelsanwendungen und Netzwerkspiele, hilfreich ist. Unter anderem wurde dazu ein neuer Algorithmus basierend auf *Circulating Envelopes* zur fälschungssicheren Bestimmung der Round-Trip-Time entwickelt.

Desweiteren wurde im Rahmen dieser Dissertation eine neuartige Architekturkomponente für Mikro-Fairness konzipiert und implementiert, die zweierlei Ziele erfüllt: zum einen erlaubt sie, Fairness für bestehende Netzwerkanwendungen ohne deren Veränderung zu erreichen und, zum anderen, dient sie als eine Forschungsplattform, mit der gezeigt wird, dass in speziellen Fällen das Fairnessproblem durch eine Transformation in ein Problem des optimalen Handels auf einem virtuellen Markt gelöst werden kann. Diese Architekturkomponente, *Fairness Proxy*, ist unser Forschungsbeitrag für einen generelleren anwendungsorientierten Ansatz für Mikro-Fairness.

Abstract

During the last decade Quality-of-Service (QoS) issues have been investigated successfully within the research community, and it only requires a further continuous process of standardization and commercial deployment to fulfill the evolving demands of modern networking applications concerning QoS.

Nonetheless, we find that the complementary issue of fairness has to date been addressed only insufficiently. With this thesis we contribute to research in that area.

We provide the link between theoretical results from the field of political economics and the needs of computer networks. Fairness definitions are given and evaluated for their use in networking. In addition, we present a first comprehensive classification of existing work on fairness in computer networks. For that reason, we divide fairness into the two dimensions: *macro-fairness* and *micro-fairness*. By *macro-fairness* we mean the distribution of network resources in a generally fair manner from which all users and applications benefit, whereas *micro-fairness* aims at fine-grained fairness for specific applications in competitive scenarios, examples of which are real-time trading systems, online auctions, etc.

Our research contribution in the field of macro-fairness is two-fold: We introduce a novel architectural proposal called *iSSD*, which achieves user-based fairness of a complexity similar to that of DiffServ, and we present our new TCP-friendly congestion control mechanism called *MPLBcc*. We validate the iSSD architecture for providing macro-fairness by means of formal analysis and network simulations. Our second research contribution to macro-fairness, MPLBcc, outperforms other TCP-

friendly mechanisms in the special case where multiple simultaneous transmission paths are available. We present both simulation results as well as field experience from the implementation of a variant of MPLBcc that was integrated as part of a bigger project for digital video transmission over the cable network return channel.

With regard to micro-fairness, we designed and validated new algorithms that achieve fair group-delay constrained traffic, which is especially useful for competitive applications such as stock trading or network games. As part of that work, we developed an algorithm for estimating round-trip time that uses *circulating envelopes* to help prevent cheating. Furthermore, we designed and implemented a novel architectural component for micro-fairness that fulfills two goals: it allows to add fairness to existing networking application scenarios without modifications to the original client-server application, and, secondly, it serves as a research platform, demonstrating that in specific cases, the problem of fairness can be solved by mapping it to the problem of optimal trading on a virtual market. This component, which we call *fairness proxy*, is our research contribution towards a more general application-layer approach to micro-fairness.

Acknowledgements

I gratefully acknowledge the support and collaboration of all the people without whom this thesis would not have been possible.

First and foremost, I wish to thank Professor Wolfgang Effelsberg for his guidance, excellent advice and many helpful suggestions throughout the years of my Ph.D. I especially thank him for his open-minded way of remotely supporting me during the final years and his ability to share a contagious fascination for the generation of innovative ideas. I also thank Professor Stefan Fischer for his suggestions during the initial phase of my thesis.

I am grateful to my colleagues Dr. Christoph Kuhmünch, Dr. Jörg Widmer and Professor Dr. Martin Mauve for fruitful discussions and our collaboration on a variety of issues related to the field of my research. Many thanks also to Betty Weyerer for her very valuable help concerning linguistic matters during the final phase of the Ph.D.

I am indebted to Dr. Albert Banchs for our brainstorming sessions on quality-of-service and fairness, during one of which he invented the terms "macro-fairness" and "micro-fairness". Albert also co-authored various papers on the initial SSD architecture. Thanks to my other former colleagues at NEC, especially to Dr. Heinrich Stüttgen, who contracted me for work related to the early phase of my Ph.D.

Furthermore, I am deeply grateful to my friend Professor Joaquín Luque for his ever warm and welcoming support, both on private and on professional matters; I also

thank all my collaborators at the University of Seville, especially for their support during the Muvitran project.

My work benefited indirectly from various sojourns in Canada and my early days of struggling with concurrent programming, for which I am indebted to Professor Peter Buhr.

Many thanks also to my friends and colleagues of the former R & D group at Enditel where I worked throughout the last three years of my Ph.D. Enditel provided me with a job that made the pursuit of my Ph.D. a really interesting challenge; without that it would have been only half the experience.

Finally, but most importantly, I want to thank my parents Monika and Ronald as well as my aunt Maria for their continuous support in all imaginable ways. Likewise, I thank Guadalupe and her family, my brother and my sister, as well as all my German, Spanish and Canadian friends for their significant non-technical help and understanding during these years. It is all of you who finally made my Ph.D. possible.

Contents

1	Introduction	1
1.1	Motivation	1
1.1.1	The Importance of Macro-Fairness	2
1.1.2	The Importance of Micro-Fairness	3
1.2	Objectives and Research Contributions	6
1.3	Outline	8
2	What is Fairness ?	9
2.1	Famous Fairness Problems	10
2.1.1	The Cake Eater Problem	10
2.1.2	The Ham Sandwich Problem	12
2.1.3	The Problem of the Nile	13
2.2	Existence of Solutions to Fairness Problems	14
2.2.1	Mathematical Background	14
2.2.2	Existence Theorems	15
2.2.2.1	Existence of Fair Divisions and Better-than-Fair Divisions	15
2.2.2.2	Existence of Best Partitions	18
2.2.3	Constructive Methods	20
2.2.3.1	Cake Division	21
2.2.3.2	Ham Sandwich Division	26

2.2.4	Discussion	30
2.3	Utility, Welfare and Fairness	30
2.3.1	Utility	31
2.3.2	Pareto Efficiency, Welfare and Fairness	32
2.3.3	Market Mechanisms	37
2.3.4	Fair Lotteries	38
2.4	Summary	38
3	Fairness Concepts in Computer Networks: A Survey	41
3.1	Introduction	42
3.1.1	The Communication Model	42
3.1.2	Utility	45
3.1.3	Examples of Network Fairness Definitions	51
3.1.3.1	Maxmin Fairness	51
3.1.3.2	Proportional Fairness	54
3.1.3.3	Utilitarian Fairness	57
3.1.3.4	The Range of Fairness Definitions: An Illustrative Example	57
3.2	Macro-Fairness - A Classification	60
3.2.1	Non-cooperative Fairness Mechanisms: Medium Access, Queue Management and Routing	63
3.2.1.1	Medium Access and Fair Scheduling	63
3.2.1.2	Fair Queue Management	68
3.2.1.3	Fair Routing	77
3.2.2	Cooperative Fairness Mechanisms: Congestion Control	78
3.2.2.1	TCP Congestion Control and Its Throughput Model	79
3.2.2.2	TCP-Friendliness	83
3.2.2.3	Fair Resource Allocation via Pricing	85
3.2.3	Fairness for Multicast	90

3.2.3.1	Congestion-controlled Multicast	90
3.2.3.2	Layered Multicast	92
3.3	Micro-Fairness - The New Challenge	93
3.3.1	Application Layer Fairness	94
3.3.2	Group-Delay Constrained Utility	96
3.3.3	Auction Mechanisms for Micro-Fairness	99
3.4	Summary	101
4	Macro-Fairness	103
4.1	iSSD: User-based Fairness for Differentiated Services	103
4.1.1	Design Considerations	104
4.1.1.1	Relative Service Differentiation	104
4.1.1.2	Stateless Core Model	105
4.1.1.3	User-based Resource Allocation	107
4.1.1.4	Real-time and Elastic Traffic Differentiation . . .	109
4.1.1.5	Ingress Control and Intra-user Discrimination . .	111
4.1.1.6	Weighted User-share-based Maxmin Fairness . .	112
4.1.1.7	Inter-ISP Boundary Adaptation	124
4.1.1.8	Compatibility to Existing Internet Mechanisms .	125
4.1.2	iSSD: The Improved Scalable Share Differentiation Architecture	126
4.1.2.1	Effective User Shares	127
4.1.2.2	Rate Estimation at Edge Nodes	128
4.1.2.3	Bandwidth Differentiation	129
4.1.2.4	Estimation of the Fair Effective Share	131
4.1.2.5	Traffic Type Separation	138
4.1.2.6	Intra-user Discrimination	140
4.1.2.7	Inter-ISP Boundaries	143
4.1.3	Simulations	145

4.1.3.1	iSSD with Different Traffic Types	146
4.1.3.2	Comparing iSSD with Other Queueing Types . .	151
4.2	MPLBcc: Multi-path Load-Balanced congestion control	154
4.2.1	Mechanisms for TCP-friendly Congestion Control	154
4.2.2	MPLBcc: Congestion Control via Multi-path Load-balancing	156
4.2.2.1	Motivation	156
4.2.2.2	Algorithm	158
4.2.2.3	Simulations	162
4.2.2.4	An Implementation Example: The MUVITRAN Framework	164
4.2.2.5	Observations and Open Issues	169
4.3	Summary	171
5	Micro-Fairness	173
5.1	Group-Delay Constrained Traffic in Client-Server Scenarios	173
5.1.1	Network Mechanisms for Fair Two-way Delay	175
5.1.2	Network Mechanisms for Fair One-way Delay	180
5.2	Fairness Proxies: Architectural Components for Adding Fairness to Existing Networks	188
5.2.1	Overview	189
5.2.2	Mechanisms	191
5.2.3	Architecture	193
5.2.4	Implementation Example: A Fairness Proxy for Fair Web Access	196
5.2.5	Measurement Results	197
5.3	Summary	205
6	Conclusions	207

Abbreviations

ABR	Available Bit Rate
ACK	Acknowledgement
AN	Access Node
APON	ATM Passive Optical Networks
ATM	Asynchronous Transfer Mode
CBQ	Class-Based Queueing
CBR	Constant Bit Rate
CSFQ	Core-Stateless Fair Queueing
CSMA/CA	Carrier Sense Multiple Access with Collision Avoidance
CTS	Clear-to-Send
CHOKe	CHOose and Keep for responsive flows, CHOose and kill for unresponsive flows
DFWMAC	Distributed Foundation Wireless Medium Access Control
DiffServ	Differentiated Services
DPS	Dynamic Packet State
DRR	Deficit Round Robin
DSLAM	Digital Subscriber Line Access Multiplexer
DWRR	Deficit Weighted Round Robin
ECN	Explicit Congestion Notification
FDM	Frequency Division Multiplexing
FEC	Foward Error Correction
FFQ	Frame-based Fair Queueing
FIFO	First-in-first-out
FLID-DL	Fair Layered Increase/Decrease with Dynamic Layering
FQ	Fair Queueing
FRED	Flow Random Early Drop
GPS	Generalized Processor Sharing
HDLC	High-level Data Link Control
IKE	Internet Key Exchange
IP	Internet Protocol
iSSD	Improved Scalable Share Differentiation
ISP	Internet Service Provider
LAN	Local Area Network
LDA+	Loss-Delay Based Adaption
LPM	Loss Path Multiplicity
LPR	Linear Proportional Response
LQD	Longest Queue Drop

LTS	Layered Transmission Scheme
MAC	Medium Access Control
MACAW	Media Access Protocol for wireless LANs
MLDA	Multicast Loss-Delay Based Adaption
MPEG	Motion Pictures Expert Group
MPLBcc	Multi-path Load-Balanced congestion control
MTCP	Multicast TCP
MUVITRAN	MUlti-path VIdeo TRANsmision
NAT	Network Address Translation
NCA	Nominee-Based Congestion Avoidance
NTP	Network Time Protocol
OBLI	Origin-Based Linear Interpolation
OLT	Optical Line Termination
ONU	Optical Network Unit
OSPF	Open Shortest Path First
PCC	Probabilistic Congestion Control
PGMCC	Pragmatic general multicast congestion control
PGPS	Packet-by-packet Generalized Processor Sharing
POTS	Plain Old Telephony Service
PQ	Priority Queueing
QoS	Quality of service
RAS	Remote Access Server
RED	Random Early Detection
RFQ	Rainbow Fair Queueing
RIO	RED with In-profile / Out-profile marking
RIP	Routing Information Protocol
RLA	Random Listening Algorithm
RLC	Receiver-Driven Layered Congestion Control
RM	Resource Management
RSVP	Resource Reservation Protocol
RTCP	Real Time Control Protocol
RTP	Real Time Protocol
RTS	Request-to-Send
RTT	Round Trip Time
SACK	Selective Acknowledgement
SBLI	Secant-Based Linear Control
SCFQ	Self-Clocked Fair Queueing
SCORE	Stateless Core
SDH	Synchronous Digital Hierarchy
SRED	Stabilized RED
TEAR	TCP Emulation at Receivers
TFRP	TCP-friendly Transport Protocol
TFRC	TCP-Friendly Rate Control Protocol
TCP	Transmission Control Protocol
TDM	Time Division Multiplexing
TUF	Tag-based Unified Fairness
UDP	User Datagram Protocol

USD	User-Share Differentiation
VC	Virtual Clock
WAN	Wide-Area Network
WDM	Wavelength Division Multiplexing
WFQ	Weighted Fair Queueing
WRR	Weighted Round Robin

List of Figures

1.1	Pyramid of users' needs in computer networks	4
2.1	A Sperner-labeled triangle	21
2.2	n -simplices for $0 \leq n \leq 3$	22
2.3	A cut-set of a cake	24
2.4	Barycentric subdivision with ownership labeling	26
2.5	Barycentric subdivision with auxiliary labeling and doors	27
2.6	Symmetric triangulation of a 3-ball with Tucker labeling	29
3.1	An abstract computer network	42
3.2	Qualitative Goodput Utility Function for elastic applications	47
3.3	Qualitative Delay Utility Function for elastic applications	48
3.4	Qualitative Delay Jitter Utility Function for elastic applications	48
3.5	Qualitative Goodput Utility Function for real-time applications	49
3.6	Qualitative Delay Utility Function for real-time applications	50
3.7	Qualitative Jitter Utility Function for real-time applications	50
3.8	A simple network scenario	53
3.9	An example of a two-node network	58
3.10	The range of fairness definitions for the two-node example	59
3.11	Comparison of Kelly's and Mo and Walrand's fairness for the two-node example	61
3.12	Classification of MAC layer schemes	64

3.13 A general switch structure	69
4.1 A core-stateless network	106
4.2 A SCORE edge node output port	107
4.3 A SCORE core node output port	107
4.4 Utility functions for a) elastic traffic and b) real-time traffic	110
4.5 Inter-ISP Boundaries	124
4.6 Traffic type separation	138
4.7 The simulation network topology	145
4.8 Throughput evolution on link 1 for simulation 1a: iSSD with real-time flows	148
4.9 Throughput evolution on link 2 for simulation 1a: iSSD with real-time flows	149
4.10 Throughput evolution on link 1 for simulation 1b: iSSD with TCP flows	150
4.11 Throughput evolution on link 2 for simulation 1b: iSSD with TCP flows	151
4.12 Multiple transmission paths	156
4.13 MPLBcc: Minimum spanning tree method	161
4.14 Throughput evolution of TFRC and MPLBcc flows on the upper bottleneck link	163
4.15 Throughput evolution of TFRC and MPLBcc flows on the lower bottleneck link	164
4.16 MUVITRAN system architecture	167
4.17 MUVITRAN laboratory set-up	169
5.1 The Circulating Envelope packet structure	177
5.2 Circulating Envelopes	178
5.3 Round-trip time calculation using Circulating Envelopes	179
5.4 Circulating Envelopes, the two-client case	180

5.5	Fair delay calculation at network nodes without per-user state	183
5.6	Queueing structure to achieve fair one-way delay	184
5.7	Algorithm summary: fair one-way delay	185
5.8	Time stamp collection for client control	186
5.9	Client-server communication without and with fairness proxy	190
5.10	High-level design of the Fairness Proxy Architecture	195
5.11	Configuration used for measurements	198
5.12	Resulting average response times without Fairness Proxy	199
5.13	Resulting average response time with Fairness Proxy, share of worst-off client: 1	200
5.14	Resulting average response time with Fairness Proxy, share of worst-off client: 15	203
5.15	Resulting average response time with Fairness Proxy, share of worst-off client: 50	204
5.16	Resulting average response time with Fairness Proxy, share of worst-off client: 75	205

List of Tables

3.1	QoS parameters for end-to-end fairness	44
3.2	Resulting bandwidth distribution in the two-node example	58
3.3	Classical auction types	100
4.1	iSSD configuration parameters for the simulations	146
4.2	Flow configuration for iSSD simulations 1a, 1b and 2	147
4.3	Results for simulation 1a: iSSD with real-time flows	149
4.4	Results for simulation 1b: iSSD with TCP flows	152
4.5	Results of simulation 2: traffic mix with different queueing mechanisms	153
4.6	Results of simulation 2: total per-user throughput	153
4.7	Resulting average throughput of TFRC and MPLBcc flows, 100 s simulation	165

