

**Why Multicast Protocols (Don't) Scale:  
An Analysis  
of Multipoint Algorithms for Scalable Group Communication**

Thesis by  
Eve M. Schooler

In Partial Fulfillment of the Requirements  
for the Degree of  
Doctor of Philosophy

Also published as Caltech Computer Science Technical Report caltechCSTR/2001.003



California Institute of Technology  
Pasadena, California

2001  
(Defended September 19, 2000)

© 2001

Eve M. Schooler

All Rights Reserved

# Acknowledgments

First and foremost, I would like to thank my committee members for their participation and feedback: Jehoshua (Shuki) Bruck, K. Mani Chandy, Deborah Estrin, Jason Hickey, and Alain Martin. I also would like to thank several other individuals who read my thesis and provided valuable comments: Bob Felderman, Joe Kiniry, Rajit Manohar, and Ruth Sivilotti. I extend special thanks to my committee chair and thesis advisor K. Mani Chandy, who has provided generous support and friendship throughout my stay at Caltech. I am also grateful to him for providing the freedom to work on whatever research presented itself as interesting. I consider myself lucky to have been part of his research group, where I came into contact with many gifted researchers: Michel Charpentier, Roman Ginis, Peter Hofstee, Joe Kiniry, K. Rustan Leino, Berna Massingill, Adam Rifkin, Paul Sivilotti, John Thornley, and Dan Zimmerman. I thank each of these individuals not only for contributing to my education, but also for enriching my experience at Caltech. Where else could one resort to the humor in using “the magic 8 ball” to make admissions or coding decisions? I have also been blessed to collaborate with Rajit Manohar, whose advice, guidance and insights have been invaluable. It is a rare friend who is willing to critique half-formed ideas.

There have been many others who keep this department running and whose help has always been greatly appreciated: Cynthia Brady, Jeri Chittum, Betta Dawson, Cindy Ferrini, Louise Foucher, Diane Goodfellow, Cici Koenig, Yvonne Recendez, and Gail Stowers. To the army of people who have maintained the systems on which I rely, I bow down to you: Dave Felt, Roman Ginis, Joe Kiniry, Dave LeBlanc, Chris Lee, Rajit Manohar, Mika Nyström, and Dan Zimmerman. To my friends at Myricom and Cornell, where most of my simulations were run, I could not have finished without the use of your fast machines (may it be a long time until I have to port the network simulator to yet another OS). I am deeply grateful. I am also thankful for the breadth of knowledge of my friends at the Fairchild Engineering Library: Kimberly Douglas and Hema Ramachandran. The resources they have assembled and have made available for on-line researchers have made my life immensely more productive.

There have been many others within the Computer Science Department who have provided both comraderie and intellectual companionship: Cindy Ball, Al Barr, Eric Bax, Cynthia Brady, David Breen, Mathieu Desbrun, Boris Dimitrov, Ilja Friedel, Eitan Grinspun, Rohit Khare, Cici

Koenig, David Laidlaw, Alain Martin, Daniel Maskit, Rajit Manohar, Mark Meyer, Mika Nyström, Marc Reiffel, Steve Taylor, Jerrell Watts, and Zöe Wood. I thank them for making this journey a memorable one. Of course, I am also especially grateful to certain unnamed individuals for never making me feel too “old.”

I owe a debt of gratitude to my office mates: Rajit Manohar for sharing his exquisite taste in chocolate and music, Zöe Wood whose artistry and spirit I prize, and Eitan Grinspun for his willingness to launch into song whenever the mood struck us!

I thank Ed Perry and his group at HP Labs Broadband Systems Lab for a rewarding summer internship in 1996. I extend heartfelt thanks to my friends and colleagues at Microsoft’s Bay Area Research Center: Jim Gemmell and Jim Gray, as well as a host of other researchers who were there during my internship during the Summer of 1997. It was my honor to work with so accomplished a group of individuals, who were willing to entertain the notion of a telecommuting internship! I have thoroughly enjoyed our collaborations.

There are many individuals from the Internet Engineering Task Force (IETF) community who I thank for providing an extra backdrop against which to do research in Networks: Mark Handley, Ruth Lang, Jörg Ott, Allison Mankin, Scott Shenker, Abel Weinrib, and many others from the Multiparty Multimedia Session Control (MMusic) working group. In addition, I extend a huge thank you to Steve Coya for making it possible for me to attend IETF meetings on a student budget.

I thank the outstanding community of people who comprise the High Speed Networking Division of USC’s Information Science Institute (ISI), where I was employed prior to Caltech and where I continued to work part-time during my first two years in graduate school. Several colleagues and friends from ISI deserve special thanks for having encouraged me over the years and for serving as constant role models: Celeste Anderson, Yigal Arens, Bob Braden, Steve Casner, Danny Cohen, Deborah Estrin, Mike Gorman, Mary Hall, and the late Jon Postel.

Other individuals whose encouragement and friendship I have valued immensely, and whose advice about staying in school I probably should have heeded back when I was completing my Masters at UCLA: Thelma and Jerry Estrin, Len and Stella Kleinrock, and Verra Morgan.

Despite appearances, there are women at Caltech! I have been fortunate to have met many phenomenal women scientists and other Caltech affiliates who have been supportive throughout: Cindy Ball, Andrea Belz, Melanie Bennett Brewer, Zehra Cataltepe, Min Chen, Christina Cohen, Roian Egnor, Carmit Eliyahu, Caroline Fohlin, Jen Linden, Berna Massingill, Helen Parker, Susan Pelletier, Cathy Wong, and Zöe Wood.

Friends I cherish for their wisdom and perspective and without whom I could not have survived: Ruth Ballenger, Suzanne Biegel, Nan Boden, Christina Cohen, Herb Donaldson, Hyewon Hyun, Rajit Manohar, Leanne Lung Nemeth, Ruth Sivilotti, Beverly Stein, and Sara Tucker.

I also wish to thank my friends who have looked after me in times of need: Dr. Lynda Roman and Nurse Charlotte Haravey of the USC Norris Cancer Center, and Nurses Divina Bautista and Alice Sogomonian at the Caltech Student Health Center.

Most importantly, I thank my family: my sisters, who have helped me to balance life and work and to understand what is truly important; my father whose intellect and curiosity know no bounds and whose fascination with invention has always inspired me; my mother who taught me my first algorithms when she taught me to read music (the for-loop of musical repeats, the go-to of skipping to a coda), who has always served as the model for the person I strive to be, and whose ability to create beauty wherever she goes awes me still.

I treasure my son Sean, who has good-naturedly marked the passage of time while I have been at Caltech! He certainly deserves his own doctorate in something: Zoology, for he can navigate the Caltech campus to find and identify most resident wildlife; Patience, for enduring my demanding schedule this past year so I could finish “the book I have been writing”; Frogology, the art of catching frogs from ponds on the Caltech campus? In addition, I could not have completed my degree were it not for Leanne Lung Nemeth, Cathy Saewert, and the wonderful teachers at the Childrens Center at Caltech, who have allowed me to do my research knowing that Sean was in good hands when he was in their care.

Above all, I am indebted to Bob Felderman, the most exceptional husband and father I know. I am grateful to him for his ceaseless encouragement, as well as his sense of humor. I thank him for happily providing advice and feedback (when solicited!), and for having the ingenuity to ask just the right questions when I was trying to debug my simulations, work out a problem, or turn a phrase. To say Bob has been infinitely supportive – not only during the normal course of everyday grad student life, but also during a particularly difficult year of health challenges – would be an understatement. We will always joke that Caltech considers spouses “dependents,” for we know who has depended on whom! For the many sacrifices he has made, and for his generosity of spirit, I am eternally grateful. It is to him that I dedicate this thesis.

The research described in this thesis was funded in part by an Earl C. Anthony Graduate Fellowship, a Career Development Grant from the American Association of University Women, a Microsoft Graduate Fellowship, as well as the Air Force Office of Scientific Research and the National Science Foundation. I thank all of them for their generous support.



# Abstract

With the exponential growth of the Internet, there is a critical need to design efficient, scalable and robust protocols to support the network infrastructure. A new class of protocols has emerged to address these challenges, and these protocols rely on a few key techniques, or micro-algorithms, to achieve scalability. By scalability, we mean the ability of groups of communicating processes to grow very large in size. We study the behavior of several of these fundamental techniques that appear in many deployed and emerging Internet standards: Suppression, Announce-Listen, and Leader Election.

These algorithms are based on the principle of efficient multipoint communication, often in combination with periodic messaging. We assume a loosely-coupled communication model, where acknowledged messaging among groups of processes is not required. Thus, processes infer information from the periodic receipt or loss of messages from other processes.

We present an analysis, validated by simulation, of the performance tradeoffs of each of these techniques. Toward this end, we derive a series of performance metrics that help us to evaluate these algorithms under lossy conditions: expected response time, network usage, memory overhead, consistency attainable, and convergence time. In addition, we study the impact of both correlated and uncorrelated loss on groups of communicating processes.

As a result, this thesis provides insights into the scalability of multicast protocols that rely upon these techniques. We provide a systematic framework for calibrating as well as predicting protocol behavior over a range of operating conditions. In the process, we establish a general methodology for the analysis of these and other scalability techniques. Finally, we explore a theory of composition; if we understand the behavior of these micro-algorithms, then we can bound analytically the performance of the more complex algorithms that rely upon them.



# Contents

<b>Acknowledgments</b>	<b>iii</b>
<b>Abstract</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Techniques for Scalability . . . . .	4
1.3 Network Model . . . . .	5
1.4 Related Work . . . . .	7
1.5 Overview . . . . .	9
<b>2 Suppression</b>	<b>11</b>
2.1 Core Algorithm . . . . .	11
2.2 Scalability . . . . .	12
2.3 Metrics . . . . .	12
2.3.1 Time Elapsed . . . . .	13
2.3.2 Extra Messages . . . . .	15
2.4 Distributions . . . . .	17
2.4.1 Uniform Distribution . . . . .	17
2.4.2 Decaying Exponential Distribution . . . . .	18
2.5 Analysis . . . . .	19
2.5.1 Realistic Parameters . . . . .	19
2.5.2 Uniform Distribution . . . . .	20
2.5.3 Decaying Exponential Distribution . . . . .	24
2.5.4 Comparisons . . . . .	24
2.6 Simulation . . . . .	28
2.7 Related Work . . . . .	30
2.8 Summary of Results . . . . .	34
2.9 Future Work . . . . .	35

<b>3</b>	<b>Suppression with Loss</b>	<b>41</b>
3.1	$E[t_{min}]$ Re-visited: Time Elapsed with Loss . . . . .	41
3.1.1	Loss Analysis . . . . .	42
3.2	Maximum Time Elapsed . . . . .	43
3.2.1	General Form . . . . .	44
3.2.2	Zero Delay . . . . .	44
3.3	Number of Messages Generated . . . . .	47
3.4	Number of Messages Required . . . . .	48
3.5	$E[\# \text{ extra}]$ Re-visited: Extra Messages with Loss . . . . .	49
3.6	Other Metrics . . . . .	51
3.7	Distributions . . . . .	51
3.7.1	Uniform Distribution . . . . .	51
3.7.2	Decaying Exponential Distribution . . . . .	52
3.8	Analysis and Simulation . . . . .	53
3.8.1	Time Metrics . . . . .	53
3.8.2	Messaging Overhead Metrics . . . . .	59
3.8.3	Loss Models . . . . .	62
3.8.4	Distributions . . . . .	62
3.9	Related Work . . . . .	66
3.10	Summary of Results . . . . .	67
3.11	Future Work . . . . .	68
<b>4</b>	<b>Announce-Listen</b>	<b>71</b>
4.1	Core Algorithm . . . . .	71
4.1.1	Scalability . . . . .	75
4.1.2	Model Parameters . . . . .	76
4.1.3	Metrics . . . . .	77
4.2	Consistency: Registry Cost . . . . .	78
4.2.1	Listener State Transition Probabilities . . . . .	79
4.2.2	Error Model . . . . .	80
4.2.3	Departures . . . . .	84
4.2.4	Overall Registry Consistency . . . . .	85
4.2.5	Arrivals . . . . .	86
4.3	Inconsistent State: Analysis and Simulation . . . . .	87
4.3.1	Analysis . . . . .	87
4.3.2	Simulation Results . . . . .	89

4.4	Convergence Time . . . . .	93
4.5	Messaging Overhead: Bandwidth . . . . .	94
4.6	Memory Overhead: Expiration Strategy . . . . .	96
4.7	Related Work . . . . .	96
4.8	Summary of Results . . . . .	98
4.9	Future Work . . . . .	99
<b>5</b>	<b>Leader Election</b>	<b>103</b>
5.1	Basic Algorithm . . . . .	104
5.1.1	The Simplest Case . . . . .	105
5.1.2	Leader Election Refined . . . . .	106
5.1.3	A Note about Timers . . . . .	107
5.2	Scalability . . . . .	108
5.3	Metrics . . . . .	108
5.3.1	Leadership Delay . . . . .	110
5.3.2	Leadership Re-establishment Delay . . . . .	112
5.3.3	Number of Messages Generated . . . . .	117
5.3.4	Inconsistent State . . . . .	122
5.4	Analysis and Simulation . . . . .	123
5.4.1	Comparison with Suppression . . . . .	124
5.4.2	Comparison with Announce-Listen . . . . .	127
5.4.3	Leader Selection Criteria . . . . .	128
5.4.4	General Trends . . . . .	130
5.5	Related Work . . . . .	131
5.6	Summary of Results . . . . .	132
5.7	Future Work . . . . .	133
<b>6</b>	<b>Conclusion</b>	<b>137</b>
<b>A</b>	<b>Announce-Listen: Inconsistency and Departures</b>	<b>141</b>
A.1	Arrivals . . . . .	141
<b>B</b>	<b>Leader Election: Rounds Needed with Uncorrelated Loss</b>	<b>145</b>
	<b>Bibliography</b>	<b>149</b>

# List of Figures

1.1	<b>Multicast vs. Unicast: Efficiency and Group Addressing.</b>	2
1.2	<b>Message Implosion Toward the Sender.</b>	3
1.3	<b>Correlated vs. Uncorrelated Loss: Message Loss Near vs. Far from Sender.</b>	6
2.1	<b>Intuition Behind <math>E[t_{min}]</math> for 2 Processes.</b>	14
2.2	<b>Suppression Algorithm.</b>	15
2.3	<b>Condition for Extra Message Transmission.</b>	15
2.4	<b>Intuition Behind <math>E[\#extra]</math> for 2 Processes.</b>	16
2.5	<b>Uniform Distribution.</b>	17
2.6	<b>Exponential Distribution.</b>	18
2.7	<b>Uniform: Time Elapsed vs. <math>T</math>.</b>	21
2.8	<b>Uniform: Time Elapsed vs. <math>N</math>.</b>	21
2.9	<b>Uniform: Extra Messages vs. <math>N</math>.</b>	22
2.10	<b>Uniform: Extra Messages vs. <math>\Delta/T</math> (Large <math>N</math>).</b>	23
2.11	<b>Uniform: Extra Messages vs. <math>\Delta/T</math> (Small <math>N</math>).</b>	23
2.12	<b>Comparison: Time Elapsed vs. <math>N</math> (Varying <math>d = \Delta</math>).</b>	24
2.13	<b>Comparison: Time Elapsed vs. <math>N</math> (Varying <math>T</math>).</b>	25
2.14	<b>Comparison: Extra Messages vs. <math>N</math> (Small <math>r = \Delta/T</math>).</b>	26
2.15	<b>Comparison: Extra Messages vs. <math>N</math> (Large <math>r = \Delta/T</math>).</b>	26
2.16	<b>Comparison: Extra Messages vs. Time Elapsed (Large <math>N</math>).</b>	27
2.17	<b>Comparison: Extra Messages vs. Time Elapsed (Small <math>N</math>).</b>	27
2.18	<b>Star Topology with fixed <math>\Delta</math>.</b>	28
2.19	<b>Simulation vs. Analysis: Time Elapsed vs. <math>T</math> (Uniform).</b>	29
2.20	<b>Simulation vs. Analysis: Time Elapsed vs. <math>N</math> (Uniform).</b>	29
2.21	<b>Simulation vs. Analysis: Extra Messages vs. <math>N</math> (Large <math>r = \Delta/T</math>).</b>	30
2.22	<b>Simulation vs. Analysis: Extra Messages vs. <math>N</math> (Small <math>r = \Delta/T</math>).</b>	30
2.23	<b>Simulation vs. Analysis: Extra Messages vs. <math>\Delta/T</math> (Exponential).</b>	31
2.24	<b>Improvement on Positive, Truncated Exponential.</b>	33
2.25	<b>Two-Component Uniform Distribution.</b>	38

3.1	<b><math>t_{max}</math>: The Maximum <math>t_{min_i}</math> Sent.</b> . . . . .	44
3.2	$Pr[t_{max} = t_{min_k}, \text{ given } i \text{ messages sent}]$ . . . . .	46
3.3	<b>Suppression Algorithm.</b> . . . . .	49
3.4	$E[t_{minr}]$ vs. $N$ ( $l = .4$ ): <b>Correlated vs. Uncorrelated.</b> . . . . .	54
3.5	$E[t_{minr}]$ vs. $N$ ( $l = .8$ ): <b>Convergence.</b> . . . . .	55
3.6	$E[t_{minr}]$ vs. $N$ ( $l = .9$ ): <b>Large Loss.</b> . . . . .	55
3.7	$E[t_{minr}]$ vs. $N$ : <b>Varying <math>l</math>.</b> . . . . .	56
3.8	$E[avg t_{max}]$ and $E[t_{minr}]$ : <b>Small Loss (<math>l = .2</math>).</b> . . . . .	56
3.9	$E[avg t_{max}]$ and $E[t_{minr}]$ : <b>Large Loss (<math>l = .95</math>).</b> . . . . .	57
3.10	<b>Probability of a Straggler in a Group of Size <math>N</math>.</b> . . . . .	58
3.11	$E[t_{max}]$ vs. $N$ : <b>Simulation vs. Analysis</b> . . . . .	59
3.12	$E[num]$ vs. $N$ ( $l = .1$ ): <b>Correlated vs. Uncorrelated.</b> . . . . .	60
3.13	$E[num]$ vs. $N$ ( $l = .7$ ): <b>Correlated vs. Uncorrelated.</b> . . . . .	60
3.14	$E[\# \text{ messages}]_{\Delta=0}$ vs. $N$ : <b>Simulation vs. Analysis.</b> . . . . .	61
3.15	$E[extra]$ vs. $N$ : <b>Correlated Loss (Varying <math>l</math>).</b> . . . . .	61
3.16	$E[t_{min k_k}]$ vs. $k$ : <b>Uniform vs. Exponential (Small <math>N</math>).</b> . . . . .	63
3.17	$E[t_{min k_k}]$ vs. $k$ : <b>Uniform vs. Exponential (Large <math>N</math>).</b> . . . . .	63
3.18	$E[t_{max}]$ vs. $N$ : <b>Correlated Loss.</b> . . . . .	64
3.19	$E[t_{max}]$ vs. $N$ : <b>Uncorrelated Loss.</b> . . . . .	64
3.20	$E[avg num]$ vs. $N$ : <b>Correlated Loss.</b> . . . . .	65
3.21	$E[required]$ vs. $N$ : <b>Uncorrelated Loss.</b> . . . . .	65
4.1	<b>Periodic Announcements and Updates.</b> . . . . .	72
4.2	<b>Announcing to Multiple Processes.</b> . . . . .	72
4.3	<b>Listener Errors.</b> . . . . .	78
4.4	<b>State Transition Probabilities.</b> . . . . .	80
4.5	<b>Event-Driven Caching and Aging Strategy.</b> . . . . .	81
4.6	<b>Boundary Condition at <math>\Delta</math>.</b> . . . . .	82
4.7	<b>Simulation vs. Analysis: Inconsistent State.</b> . . . . .	91
4.8	<b>Simulation: Inconsistent State (Large ratio = <math>\Delta/T</math>).</b> . . . . .	91
4.9	<b>Simulation: Inconsistent State (Small ratio = <math>\Delta/T</math>).</b> . . . . .	91
4.10	<b>Simulation: Data Change Rate (<math>\Delta/T = .1</math>).</b> . . . . .	92
4.11	<b>Simulation: Data Change Rate (<math>\Delta/T = .01</math>).</b> . . . . .	92
4.12	<b>Simulation: Data Change Rate (<math>\Delta/T = .001</math>).</b> . . . . .	92
4.13	<b>Convergence Time.</b> . . . . .	94
4.14	<b>Roundtrip Time Estimate: Calculating Pair-wise <math>\Delta</math>.</b> . . . . .	102

5.1	<b>Announce Leadership.</b>	106
5.2	<b>Suppress Leadership.</b>	108
5.3	<b>Time to Notice the Leader Left (No Loss).</b>	113
5.4	<b>Leadership Delay with a Late Joiner: <math>joiner \neq leader</math></b>	114
5.5	<b>Vector <math>\vec{Q}</math> of Process Addresses Ordered from Largest to Smallest.</b>	118
5.6	<b>LE vs. SUP: <math>E[t_{max}]</math> vs. <math>N</math> (<math>l = 0</math>).</b>	124
5.7	<b>LE vs. SUP: <math>E[num\ msgs]</math> vs. <math>N</math> (<math>l = 0</math>).</b>	125
5.8	<b>LE vs. SUP: <math>E[t_{max}]</math> vs. <math>N</math> (<math>l = .4</math>).</b>	126
5.9	<b>LE vs. SUP: <math>E[num\ msgs]</math> vs. <math>N</math> (<math>l = .4</math>).</b>	126
5.10	<b><math>E[t_{max}]</math> vs. <math>N</math>: Correlated Loss (Varying <math>l</math>).</b>	127
5.11	<b>LE vs. AL: <math>E[t_{max}]</math> vs. <math>N</math> (<math>l = .4</math>).</b>	128
5.12	<b>First Round Comparison of LE-M with LE and SUP: <math>E[t_{max}]</math> vs. <math>N</math> (<math>l = .4</math>).</b>	129
5.13	<b>LE-M Compared with LE and SUP: <math>E[t_{max}]</math> vs. <math>N</math> (<math>l = .4</math>).</b>	130

# List of Tables

2.1	<b>Timer Range (<math>T</math>).</b> . . . . .	20
2.2	<b>Number of Processes (<math>N</math>).</b> . . . . .	20
2.3	<b>Transmission Delay (<math>\Delta</math>).</b> . . . . .	20
3.1	<b>Zero-Delay Event Probabilities.</b> . . . . .	45
3.2	<b>Time Metrics Simulated.</b> . . . . .	53
3.3	<b>Messaging Overhead Metrics Simulated.</b> . . . . .	59
4.1	<b>Announce-Listen Parameters.</b> . . . . .	77
4.2	<b>Inconsistent State: <math>nT \leq t \leq nT + \Delta</math>.</b> . . . . .	83
4.3	<b>Inconsistent State: <math>nT + \Delta \leq t \leq (n + 1)T</math>.</b> . . . . .	83
4.4	<b>Announcement Message Content.</b> . . . . .	90
5.1	<b>Convergence Metrics Simulated.</b> . . . . .	124
A.1	<b>Inconsistency: After <math>\Delta</math>.</b> . . . . .	142
A.2	<b>Inconsistency: Within <math>\Delta</math>.</b> . . . . .	142

# List of Programs

2.1	Suppression Algorithm. . . . .	12
3.1	Parameterized Suppression Algorithm. . . . .	69
4.1	Announce-Listen Algorithm. . . . .	73
4.2	Announce-Listen Algorithm with Caching. . . . .	73
5.1	Leader Election Algorithm using Announce-Listen. . . . .	105
5.2	Leader Election Algorithm using Suppression. . . . .	107