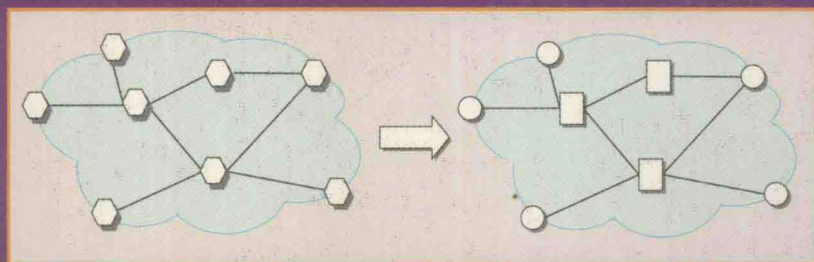


Ion Stoica

# Stateless Core: A Scalable Approach for Quality of Service in the Internet

Winning Thesis of the  
2001 ACM Doctoral Dissertation Competition



Springer

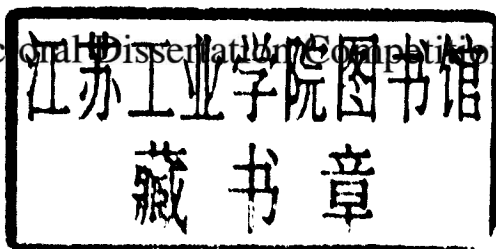


ASSOCIATION FOR  
COMPUTING MACHINERY

Ion Stoica

# Stateless Core: A Scalable Approach for Quality of Service in the Internet

Winning Thesis of  
the 2001 ACM Doctoral Dissertation Competition



Springer

## Series Editors

Gerhard Goos, Karlsruhe University, Germany  
Juris Hartmanis, Cornell University, NY, USA  
Jan van Leeuwen, Utrecht University, The Netherlands

## Author

Ion Stoica  
University of California at Berkeley, Computer Science Division  
645 Soda Hall, Berkeley, CA 94720-1776, USA  
E-mail: istoica@cs.berkeley.edu

Library of Congress Control Number: 2004104504

CR Subject Classification (1998): C.2, E.1, H.3.4-5, H.4.3

ISSN 0302-9743

ISBN 3-540-21960-9 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable to prosecution under the German Copyright Law.

Springer-Verlag is a part of Springer Science+Business Media

springeronline.com

© Springer-Verlag Berlin Heidelberg 2004  
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Boller Mediendesign  
Printed on acid-free paper SPIN: 11007180 06/3142 5 4 3 2 1 0

# Foreword

The fundamental aspect of the Internet architecture that distinguishes it from other network technologies (such as X.25 and ATM) is that it is connectionless (vs. connection-oriented) and stateless (vs. stateful). The heated debate of whether connection-oriented or connectionless architecture is better has lasted for several decades. Proponents of the connectionless architecture point out the great robustness and scalability properties of the architecture, as demonstrated by the Internet. One well-known articulation of this philosophy is the “End-to-End Arguments”. Opponents argue, rightfully, that there is no known solution that can provide quantitative performance assurances or guaranteed QoS in a connectionless network. It has been widely recognized that QoS is a must-have feature as the Internet technology evolves to the next stage. However, all existing solutions that provide guaranteed QoS require routers to maintain per flow (another name for connection used by the Internet community) state, which is the fundamental element of a connection-oriented architecture. The apparent conflicting goals of having a stateless network and supporting QoS have presented a great dilemma for Internet architects. As an example, Dave Clark, one of the most respected Internet architects and the author of the famous “End-to-End Arguments” paper, was also a key designer of the Internet Integrated Services Architecture that requires routers to maintain per flow state.

Dr. Ion Stoica’s dissertation addresses this most pressing and difficult problem facing the Internet community today: how to enhance the Internet to support rich functionalities (such as QoS and traffic management) while still maintaining the scalability and robustness properties embodied in the original Internet architecture.

In his dissertation, Dr. Stoica proposes a novel architecture called SCORE (Stateless Core) that does not require core routers to maintain per flow state yet can provide services similar to those provided by stateful networks. This is achieved by a family of SCORE distributed algorithms that approximate the services provided by idealized stateful networks. The key technique used to implement a SCORE network is Dynamic Packet State (DPS), which uses extra state carried in packet headers to coordinate distributed algorithms implemented by routers. Such an architecture has both important theoretical and practical significances. From a conceptual point of view, this architecture

## VIII Foreword

is the first that combines the advantages of stateful and stateless networks, and can therefore achieve QoS, scalability, and robustness simultaneously. From a practical point of view, the industry and the IETF have been struggling to make a choice between two QoS architectures: the stateful Intserv, which can provide hard QoS guarantees but is less scalable and robust, and the stateless Diffserv, which is more scalable and robust but cannot provide services with high assurance. The SCORE architecture provides a third approach that is superior.

I believe that this research represents one of the most important and innovative contributions in networking research in the past decade. I hope that you will enjoy reading it and agree with me afterward.

*Hui Zhang*

# Preface

This book contains the dissertation the author wrote at the Department of Electrical and Computer Engineering (ECE) at Carnegie Mellon University. This thesis was submitted to the ECE department in conformity with the requirements for the degree of Doctor of Philosophy in 2000. It was honored with the 2001 ACM Doctoral Dissertation Award.

## Abstract

Today's Internet provides one simple service: best-effort datagram delivery. This minimalist service allows the Internet to be *stateless*, that is, routers do not need to maintain any fine-grained information about traffic. As a result of this stateless architecture, the Internet is both highly *scalable* and *robust*. However, as the Internet evolves into a global commercial infrastructure that is expected to support a plethora of new applications, such as IP telephony, interactive TV, and e-commerce, the existing best-effort service will no longer be sufficient. As a consequence, there is an urgent need to provide more powerful services such as guaranteed services, differentiated services, and flow protection.

Over the past decade, there has been intense research toward achieving this goal. Two classes of solutions have been proposed: those maintaining the *stateless* property of the original Internet (e.g., differentiated services), and those requiring a new *stateful* architecture (e.g., integrated services). While stateful solutions can provide more powerful and flexible services, such as per flow bandwidth and delay guarantees, they are less scalable than stateless solutions. In particular, stateful solutions require each router to maintain and manage per flow state on the control path, and to perform per flow classification, scheduling, and buffer management on the data path. Since today's routers can handle millions of active flows, it is difficult, if not impossible, to implement such solutions in a scalable fashion. On the other hand, while stateless solutions are much more scalable, they offer weaker services.

The key contribution of this dissertation is to bridge this long-standing gap between stateless and stateful solutions in packet-switched networks such as the Internet. Our thesis is that "*it is actually possible to provide services as powerful and as flexible as the ones implemented by a stateful network using a*



*stateless network*". To prove this thesis, we propose a novel technique called Dynamic Packet State (DPS). The key idea behind DPS is that, instead of having routers maintain per flow state, packets carry the state. In this way, routers are still able to process packets on a per flow basis, despite the fact that they do not maintain any per flow state. Based on DPS, we develop a network architecture called Stateless Core (SCORE) in which core routers do not maintain any per flow state. Yet, using DPS to coordinate actions of edge and core routers along the path traversed by a flow allows us to design distributed algorithms that emulate the behavior of a broad class of stateful networks in SCORE networks.

In this dissertation we describe complete solutions including architectures, algorithms, and implementations which address three of the most important problems in today's Internet: providing guaranteed services, differentiated services, and flow protection. Compared to existing solutions, our solutions *eliminate* the most complex operations on both the data and control paths in the network core, i.e., packet classification on the data path, and maintaining per flow state consistency on the control path. In addition, the complexities of buffer management and packet scheduling are greatly reduced. For example, in our flow protection solution these operations take constant time, while in previous solutions these operations may take time logarithmic in the number of flows traversing the router.

## Acknowledgements

I am extremely grateful to Hui Zhang, my thesis advisor, for giving me the right amount of freedom and guidance during my graduate studies. From the very beginning, he treated me as a peer and as a friend. He was instrumental in maintaining my focus, and constantly reminding me that identifying the research problem is as important, if not more important, than finding the right solution. Hui not only taught me how to become a better researcher, but also helped me to become a better person. His engaging arguments and strong feedback contributed greatly to this dissertation. I hope for and look forward to continued collaboration with him in the future.

The genesis of this thesis can be traced back to my internship at Xerox PARC in the summer of 1997. It all started with Scott Shenker asking the intriguing question: "Can we approximate Fair Queueing without maintaining per flow state in a network cloud?" I am indebted to Scott for teaching me how to rigorously define a problem and then pursue its solution. During these years he was an invaluable source of feedback and support. His suggestions and insights had a great impact on this dissertation.

I am grateful to the other members of my committee, Garth Gibson, Thomas Gross, and Peter Steenkiste, for their feedback and for their advice that helped to shape my research skills. Garth taught me the art of asking the right questions in a technical discourse. His inquisitorial and sharp questions

helped me to better understand the limitations of my research and motivated me to find better ways to explain my results. Thomas provided the right balance to my research by constantly encouraging me to not get buried in the algorithmic details, but to always try to put my work into perspective. Peter always found time to discuss research issues, and gave excellent feedback. He was one of the first to suggest using Dynamic Packet State to provide guaranteed services.

Thanks to Mihai Budiu, Yang-hua Chu, Urs Hengartner, Eugene Ng, Mahadev Satyanarayanan, and Jeannette Wing for their feedback and comments that helped to improve the overall quality of this dissertation. Thanks to Joan Digney for accommodating me in her busy schedule and proofreading this thesis, which helped to significantly improve the presentation.

I am grateful to all the friends and colleagues with whom I spent my time as a graduate student at Carnegie Mellon University. Thanks to my Romanian friends, Mihai and Raluca Budiu, Cristian Dima, Marius Minea, and George Necula, with whom I spent many hours discussing the most varied and exciting topics. Mihai sparkled many of these conversations with his wit, and by being a never empty reservoir of information. Thanks to my networking group colleagues and officemates Yang-hua Chu, Urs Hengartner, Nick Hopper, Tom Kang, Marco Mellia, Andrew Myers, Eugene Ng, Sanjay Rao, Chuck Rosenberg, Kay Sripanidkulchai, Donpaul Stephens, and Yinglian Xie. I treasure our animated Friday lunch discussions that always managed to end the week on a high note. Special thanks to Eugene for his help and patience with my countless questions. He was the default good-to-answer-all-questions person whom I asked about everything, from how to modify the `if_de` driver in FreeBSD, to which are the best restaurants around the campus.

The completion of this thesis marked the end of my many years as a student. Among many outstanding teachers who shaped my education and scientific career are: Hussein Abdel-Wahab, Irina Athanasiu, Kevin Jeffay, David Keyes, Trandafir Moisa, Stephan Olariu, Alex Pothen, and Nicolae Tapus. I am grateful to Kevin, who was instrumental in helping and then convincing me to come to Carnegie Mellon University. Many thanks to Irina who, during my studies at the “Politehnica” University of Bucharest, gave me the mentorship a student can only dream about.

I would like to express my earnest gratitude to my parents and my sister for their love and support, without which any of my achievements would not have been possible. Thanks to my father who sparked my early interest in science and engineering. His undaunting confidence gave me the strength to overcome any difficulties and to maintain high goals. Thanks to my mother for her love and countless sacrifices to raise me and to give me the best possible education.

I am deeply indebted to my dear wife Emilia for her love and understanding through my graduate years. She was always behind me and gave her unconditional support even if that meant sacrificing the time we spent



## XII Preface

together. I thank my mother and my mother-in-law who devotedly took care of our son for two and a half years. Without their help, it would not have been possible to reach this stage in my career. Finally, thanks to our son, George, for the joy and the happiness he brings to me during our many moments together.

Berkeley, December 2003

*Ion Stoica*

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Main Contribution	2
1.2	Other Contributions	5
1.3	Evaluation	6
1.4	Discussion	7
1.4.1	Why Per Flow Processing?	7
1.4.2	Scalability Concerns with Stateful Network Architectures	9
1.5	Organization	10
<b>2</b>	<b>Background</b>	<b>13</b>
2.1	Circuit Switching Vs. Packet Switching	14
2.2	IP Network Model	15
2.2.1	Router Architecture	15
2.2.2	Data Path	17
2.2.3	Control Path	21
2.2.4	Discussion	24
2.3	Network Service Taxonomy	25
2.3.1	Best Effort Service	26
2.3.2	Flow Protection: Network Support for Congestion Control	27
2.3.3	Integrated Services	29
2.3.4	Differentiated Services	30
2.4	Summary	33
<b>3</b>	<b>Overview</b>	<b>35</b>
3.1	Solution Overview	35
3.1.1	The Stateless Core (SCORE) Network Architecture	36
3.1.2	The “State-Elimination” Approach	36
3.1.3	The Dynamic Packet State (DPS) Technique	36
3.2	Prototype Implementation	43
3.2.1	An Example	44

3.3	Comparison to Intserv and Diffserv .....	47
3.3.1	Intserv .....	47
3.3.2	Diffserv .....	49
3.4	Summary .....	51
<b>4</b>	<b>Providing Flow Protection in SCORE .....</b>	<b>53</b>
4.1	Background .....	53
4.2	Solution Outline .....	55
4.3	Core-Stateless Fair Queueing (CSFQ) .....	56
4.3.1	Fluid Model Algorithm .....	56
4.3.2	Packet Algorithm .....	57
4.3.3	Weighted CSFQ .....	61
4.3.4	Performance Bounds .....	61
4.3.5	Implementation Complexity .....	62
4.3.6	Architectural Considerations .....	62
4.3.7	Miscellaneous Details .....	63
4.4	Simulation Results .....	63
4.4.1	A Single Congested Link .....	65
4.4.2	Multiple Congested Links .....	66
4.4.3	Coexistence of Different Adaptation Schemes .....	67
4.4.4	Different Traffic Models .....	70
4.4.5	Large Latency .....	71
4.4.6	Packet Relabeling .....	72
4.4.7	Discussion of Simulation Results .....	73
4.5	Related Work .....	73
4.6	Summary .....	75
<b>5</b>	<b>Providing Guaranteed Services in SCORE .....</b>	<b>77</b>
5.1	Background .....	77
5.2	Solution Outline .....	79
5.3	Data Plane: Scheduling without Per Flow State .....	80
5.3.1	Jitter Virtual Clock (Jitter-VC) .....	80
5.3.2	Core-Jitter-VC (CJVC) .....	82
5.3.3	Data Path Complexity .....	85
5.4	Control Plane: Admission Control with no Per Flow State ...	87
5.4.1	Ingress-to-Egress Admission Control .....	88
5.4.2	Per-Hop Admission Control .....	89
5.4.3	Aggregate Reservation Estimation Algorithm .....	90
5.5	Experimental Results .....	95
5.5.1	Processing Overhead .....	100
5.6	Related Work .....	100
5.7	Summary .....	102

<b>6</b>	<b>Providing Relative Service Differentiation in SCORE . . . .</b>	<b>103</b>
6.1	Background . . . . .	104
6.2	Solution Outline . . . . .	106
6.3	LIRA: Service Differentiation Based on Resource Right Tokens . . . . .	107
6.3.1	Link Cost Computation . . . . .	109
6.3.2	Path Cost Computation and Distribution . . . . .	110
6.3.3	Multipath Routing and Load Balancing . . . . .	111
6.3.4	Route Pinning . . . . .	111
6.3.5	Path Selection . . . . .	113
6.3.6	Scalability . . . . .	113
6.4	Simulation Results . . . . .	114
6.4.1	Experiment Design . . . . .	115
6.4.2	Experiment 1: Local Fairness and Service Differentiation . . . . .	117
6.4.3	Experiment 2: User Fairness and Load Balancing . . . .	119
6.4.4	Experiment 3: Load Distribution and Load Balancing .	119
6.4.5	Experiment 4: Large Scale Example . . . . .	121
6.4.6	Summary of Simulation Results . . . . .	123
6.5	Discussion . . . . .	124
6.6	Related Work . . . . .	125
6.7	Summary . . . . .	126
<b>7</b>	<b>Making SCORE More Robust and Scalable . . . . .</b>	<b>129</b>
7.1	Failure Model . . . . .	130
7.1.1	Example . . . . .	131
7.2	The “Verify-and-Protect” Approach . . . . .	132
7.2.1	Node Identification . . . . .	133
7.2.2	Protection . . . . .	133
7.2.3	Recovery . . . . .	134
7.3	Flow Verification . . . . .	134
7.3.1	Bufferless Packet System . . . . .	137
7.3.2	Flow Identification Test . . . . .	137
7.3.3	Setting Threshold $H_u$ . . . . .	140
7.3.4	Increasing Flow Identification Test’s Robustness and Responsiveness . . . . .	142
7.4	Identifying Misbehaving Nodes . . . . .	143
7.4.1	General Properties . . . . .	144
7.5	Simulation Results . . . . .	146
7.5.1	Calibration . . . . .	147
7.5.2	Protection and Recovery . . . . .	148
7.6	Summary . . . . .	150

<b>8</b>	<b>Prototype Implementation Description</b>	153
8.1	Prototype Implementation	154
8.1.1	Updating State in IP Header	155
8.1.2	Data Path	156
8.1.3	Control Path	158
8.2	Carrying State in Data Packets	158
8.2.1	Carrying State in IP Header	160
8.2.2	Efficient State Encoding	160
8.2.3	State Encoding for Guaranteed Service	162
8.2.4	State Encoding for LIRA	165
8.2.5	State Encoding for CSFQ	166
8.2.6	State Encoding Formats for Future Use	166
8.3	System Monitoring	166
8.4	System Configuration	168
8.4.1	Router Configuration	169
8.4.2	Flow Reservation	169
8.4.3	Monitoring	170
8.5	Summary	170
<b>9</b>	<b>Conclusions and Future Work</b>	173
9.1	Contributions	173
9.2	Limitations	175
9.3	Future Work	177
9.3.1	Decoupling Bandwidth and Delay Allocations	177
9.3.2	Excess Bandwidth Allocation	178
9.3.3	Link Sharing	179
9.3.4	Multicast	181
9.3.5	Verifiable End-to-End Protocols	181
9.3.6	Incremental Deployability	182
9.3.7	General Framework	183
9.4	Final Remarks	183
<b>A</b>	<b>Performance Bounds for CSFQ</b>	185
<b>B</b>	<b>Performance Bounds for Guaranteed Services</b>	191
B.1	Network Utilization of Premium Service in Diffserv Networks	191
B.2	Proof of Theorem 2	193
B.3	Proof of Theorem 3	198
B.3.1	Identical Flow Rates	198
B.3.2	Arbitrary Flow Rates	204
B.4	Proof of Theorem 4	211
	<b>References</b>	213

# 1 Introduction

Today's Internet provides one simple service: best effort datagram delivery. Such a minimalist service allows routers to be *stateless*, that is, except for the routing state, which is highly aggregated, routers do not need to maintain any fine grained state about traffic. As a consequence, today's Internet is both highly *scalable* and *robust*. It is scalable because router complexity does not increase in either the number of flows or the number of nodes in the network, and it is robust because there is little state, if any, to update when a router fails or recovers. The scalability and robustness are two of the most important reasons behind the success of today's Internet.

However, as the Internet evolves into a global commercial infrastructure, there is a growing need to provide more powerful services than best effort such as guaranteed services, differentiated services, and flow protection. Guaranteed services would make it possible to guarantee performance parameters such as bandwidth and delay on a per flow basis. An example would be to guarantee that a flow receives at least a specified amount of bandwidth, ensuring that the delay experienced by its packets does not exceed a specified threshold. This service would provide support for new applications such as IP telephony, video-conferencing, and remote diagnostics. Differentiated services would allow us to provide bandwidth and loss rate differentiation for traffic aggregates over multiple granularities ranging from individual flows to the entire traffic of a large organization. An example would be to allocate to one organization twice as much bandwidth on every link in the network as another organization. Flow protection would allow diverse end-to-end congestion control schemes to seamlessly coexist in the Internet, protecting the well behaved traffic from the malicious or ill-behaved traffic. For example, if two flows share the same link, with flow protection, each flow will get at least half of the link capacity independent of the behavior of the other flow, as long as the flow has enough demand. In contrast, in today's Internet, a malicious flow that sends traffic at a higher rate than the link capacity can provoke packet losses to another flow no matter how little traffic that flow sends!

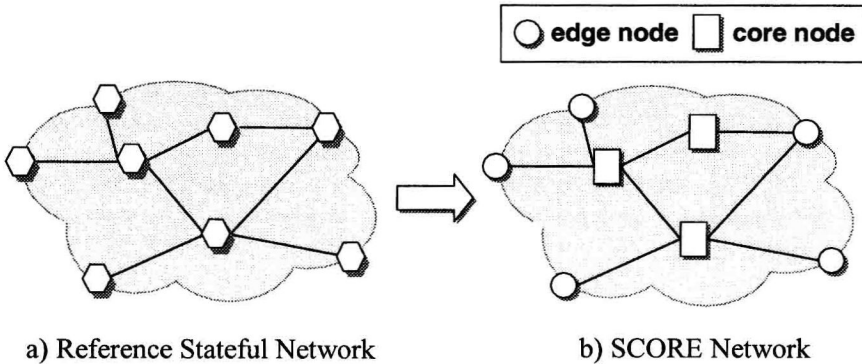
Providing these services in packet switched networks such as the Internet has been one of the major challenges in the network research over the past decade. To address this challenge, a plethora of techniques and mechanisms



have been developed for packet scheduling, buffer management, and signaling. While the proposed solutions are able to provide very powerful network services, they come at a cost: *complexity*. In particular, these solutions usually assume a *stateful* network architecture, that is, a network in which every router maintains per flow state. Since there can be a large number of active flows in the Internet, and this number is expected to continue to increase at an exponential rate, it is an open question whether such an architecture can be efficiently implemented. In addition, due to the complex algorithms required to set and preserve the state consistency across the network, robustness is much harder to achieve.

In summary, while stateful architectures can provide more sophisticated services than the best effort service, stateless architectures such as the current Internet are more scalable and robust. The natural question is then: Can we achieve the best of the two worlds? That is, *is it possible to provide services as powerful and flexible as the ones implemented by a stateful network in a stateless network?*

In this dissertation we answer this question affirmatively by showing that some of the most representative Internet services that require stateful networks can indeed be implemented in a mostly stateless network architecture.



**Fig. 1.1.** (a) A reference stateful network whose functionality is approximated by (b) a Stateless Core (SCORE) network. In SCORE only edge nodes maintain per flow state and perform per flow management; core nodes do not maintain any per flow state.

## 1.1 Main Contribution

The main contribution of this dissertation is to *provide the first solution that bridges the long-standing gap between stateless and stateful network architectures*. In particular, we show that three of the most important Internet

services proposed in literature during the past decade, and for which the previous known solutions require stateful networks, can be implemented in a stateless core network. These services are: (1) guaranteed services, (2) service differentiation for large granularity traffic, and (3) flow protection to provide network support for congestion control.

The main goal of our solution is to push the state and therefore the complexity out of the network core, *without* compromising network ability to provide per flow services. The key ideas that allow us to achieve this goal are:

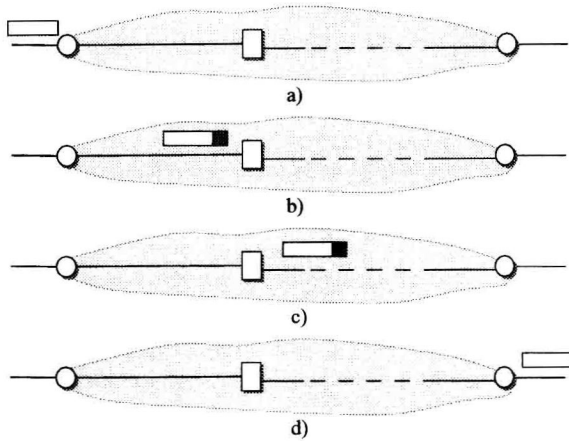
1. instead of having core nodes maintain per flow state, have packets carry this state, and
2. use the state carried by the packets to implement distributed algorithms to provide network services as powerful and as flexible as the ones implemented by stateful networks

The following paragraphs present the main components of our solution:

**The Stateless Core (SCORE) Network Architecture** The basic building block of our solution is the Stateless Core (SCORE) domain. We define a SCORE domain as being a trusted and contiguous region of network in which only edge routers maintain per flow state; the core routers do *not* maintain any per flow state (see Figure 1.1(b)). Since edge routers usually run at a much lower speed and handle far fewer flows than core routers, this architecture is highly scalable.

**The “State-Elimination” Approach** Our ultimate goal is to provide powerful and flexible network services in a stateless network architecture. To achieve this goal, we propose an approach, called “state-elimination” approach, that consists of two steps (see Figure 1.1). The first step is to define a reference stateful network that implements the desired service. The second step is to approximate or, if possible, to emulate the functionality of the reference network in a SCORE network. By doing this, we can provide services as powerful and flexible as the ones implemented by a stateful network in a mostly stateless network architecture, i.e., in a SCORE network.

**The Dynamic Packet State (DPS) Technique** To implement the approach, we propose a novel technique called Dynamic Packet State (DPS). As shown in Figure 1.2, with DPS, each packet carries in its header some state that is initialized by the ingress router. Core routers process each incoming packet based on the state carried in the packet’s header, updating both its internal state and the state in the packet’s header before forwarding it to the next hop. In this way, routers are able to process packets on a per flow basis, despite the fact that they do not maintain per flow state. As we will demonstrate in this dissertation, by using DPS to coordinate the actions of edge and core routers along the path traversed by a flow, it is possible to



**Fig. 1.2.** An illustration of the Dynamic Packet State (DPS) technique used to implement per flow services in a SCORE network: (a-b) upon a packet arrival the ingress node inserts some flow dependent state into the packet header; (b-c) a core node processes the packet based on this state, and eventually updates both its internal state and the packet state before forwarding it. (c-d) the egress node removes the state from the packet header.

design distributed algorithms to approximate the behavior of a broad class of stateful networks using networks in which core routers do not maintain per flow state.

**The “Verify-and-Protect” Approach** While our solutions based on SCORE/DPS have many advantages over traditional stateful solutions, they still suffer from robustness and scalability limitations when compared to stateless solutions. The scalability of the SCORE architecture suffers from the fact that the network core cannot transcend trust boundaries (such as boundaries between competing Internet Service Providers), and therefore high-speed routers on these boundaries must be stateful edge routers. System robustness is limited by the possibility that a single edge or core router may malfunction, inserting erroneous information in the packet headers, thus severely impacting performance of the entire network.

In Chapter 7 we propose an approach, called “verify-and-protect”, that overcomes these limitations. We achieve scalability by pushing the complexity all the way to the end-hosts, eliminating the distinction between edge and core routers. To address the trust and robustness issues, all routers statistically verify that the incoming packets are correctly marked. This approach enables routers to discover and isolate misbehaving end-hosts and routers.