

2ND  
EDITION

# THE BOOK OF PF

A NO-NONSENSE GUIDE TO THE  
OPENBSD FIREWALL

PETER N.M. HANSTEEN





## **PRAISE FOR THE FIRST EDITION OF *THE BOOK OF PF***

“This book is for everyone who uses PF. Regardless of operating system and skill level, this book will teach you something new and interesting.”

—BSD MAGAZINE

“With Mr. Hansteen paying close attention to important topics like state inspection, SPAM, black/grey listing, and many others, this must-have reference for BSD users can go a long way to helping you fine tune the who/what/where/when/how of access control on your BSD box.”

—INFOWORLD

“A must-have resource for anyone who deals with firewall configurations. If you’ve heard good things about PF and have been thinking of giving it a go, this book is definitely for you. Start at the beginning and before you know it you’ll be through the book and quite the PF guru. Even if you’re already a PF guru, this is still a good book to keep on the shelf to refer to in thorny situations or to lend to colleagues.”

—DRU LAVIGNE, TECH WRITER

“The book is a great resource and has me eager to rewrite my aging rulesets.”

—;LOGIN:

“This book is a super-easy read. I loved it! This book easily makes my Top 5 Book list.”

—DAEMON NEWS



# **THE BOOK OF™ PF**

## **2ND EDITION**

**A NO-NONSENSE GUIDE TO THE  
OPENBSD FIREWALL**

**by Peter N.M. Hansteen**



San Francisco

**THE BOOK OF PF, 2ND EDITION.** Copyright © 2011 by Peter N.M. Hansteen.

All rights reserved. No part of this work may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval system, without the prior written permission of the copyright owner and the publisher.

14 13 12 11 10    1 2 3 4 5 6 7 8 9

ISBN-10: 1-59327-274-X

ISBN-13: 978-1-59327-274-6

Publisher: William Pollock

Production Editors: Ansel Staton and Serena Yang

Cover and Interior Design: Octopod Studios

Developmental Editor: William Pollock

Technical Reviewer: Henning Brauer

Copyeditor: Marilyn Smith

Compositors: Riley Hoffman and Ansel Staton

Proofreader: Linda Seifert

Indexer: Valerie Haynes Perry

For information on book distributors or translations, please contact No Starch Press, Inc. directly:

No Starch Press, Inc.

38 Ringold Street, San Francisco, CA 94103

phone: 415.863.9900; fax: 415.863.9950; info@nostarch.com; www.nostarch.com

*The Library of Congress has cataloged the first edition as follows:*

Hansteen, Peter N. M.

The book of PF : a no-nonsense guide to the OpenBSD firewall / Peter N.M. Hansteen.

p. cm.

Includes index.

ISBN-13: 978-1-59327-165-7

ISBN-10: 1-59327-165-4

1. OpenBSD (Electronic resource) 2. TCP/IP (Computer network protocol) 3. Firewalls (Computer security)

I. Title.

TK5105.585.H385 2008

005.8--dc22

2007042929

No Starch Press and the No Starch Press logo are registered trademarks of No Starch Press, Inc. Other product and company names mentioned herein may be the trademarks of their respective owners. Rather than use a trademark symbol with every occurrence of a trademarked name, we are using the names only in an editorial fashion and to the benefit of the trademark owner, with no intention of infringement of the trademark.

The information in this book is distributed on an “As Is” basis, without warranty. While every precaution has been taken in the preparation of this work, neither the author nor No Starch Press, Inc. shall have any liability to any person or entity with respect to any loss or damage caused or alleged to be caused directly or indirectly by the information contained in it.

To Gene Scharmann,  
who all those years ago nudged me in the direction of free software



## BRIEF CONTENTS

Foreword by Bob Beck (from the first edition) .....	xiii
Acknowledgments .....	xv
Introduction .....	xvii
Chapter 1: Building the Network You Need .....	1
Chapter 2: PF Configuration Basics .....	11
Chapter 3: Into the Real World .....	25
Chapter 4: Wireless Networks Made Easy .....	41
Chapter 5: Bigger or Trickier Networks.....	59
Chapter 6: Turning the Tables for Proactive Defense .....	85
Chapter 7: Queues, Shaping, and Redundancy .....	105
Chapter 8: Logging, Monitoring, and Statistics .....	131
Chapter 9: Getting Your Setup Just Right.....	151
Appendix A: Resources .....	167
Appendix B: A Note on Hardware Support .....	173
Index .....	177



# CONTENTS IN DETAIL

**FOREWORD by Bob Beck (from the first edition)** **xiii**

**ACKNOWLEDGMENTS** **xv**

**INTRODUCTION** **xvii**

This Is Not a HOWTO .....	xviii
What This Book Covers .....	xviii

**1** **BUILDING THE NETWORK YOU NEED** **1**

Your Network: High Performance, Low Maintenance, and Secure .....	1
Where the Packet Filter Fits In .....	3
The Rise of PF .....	3
If You Came from Elsewhere .....	5
Pointers for Linux Users .....	6
Frequently Answered Questions About PF .....	7
A Little Encouragement: A PF Haiku .....	9

**2** **PF CONFIGURATION BASICS** **11**

The First Step: Enabling PF .....	12
Setting Up PF on OpenBSD .....	12
Setting Up PF on FreeBSD .....	13
Setting Up PF on NetBSD .....	15
A Simple PF Rule Set: A Single, Stand-Alone Machine .....	16
A Minimal Rule Set .....	16
Testing the Rule Set .....	18
Slightly Stricter: Using Lists and Macros for Readability .....	18
A Stricter Baseline Rule Set .....	19
Reloading the Rule Set and Looking for Errors .....	20
Checking Your Rules .....	21
Testing the Changed Rule Set .....	21
Displaying Information About Your System .....	22
Looking Ahead .....	23

**3** **INTO THE REAL WORLD** **25**

A Simple Gateway .....	26
Keep It Simple: Avoid the Pitfalls of in, out, and on .....	26
Network Address Translation vs. IPv6 .....	27

Final Preparations: Defining Your Local Network .....	28
Setting Up a Gateway .....	29
Testing Your Rule Set .....	33
That Sad Old FTP Thing .....	34
If We Must: <i>ftp-proxy</i> with Redirection .....	34
Making Your Network Troubleshooting Friendly .....	36
Do We Let It All Through? .....	37
The Easy Way Out: The Buck Stops Here .....	37
Letting ping Through .....	37
Helping traceroute .....	38
Path MTU Discovery .....	38
Tables Make Your Life Easier .....	39

## 4 WIRELESS NETWORKS MADE EASY 41

A Little IEEE 802.11 Background .....	42
MAC Address Filtering .....	42
WEP .....	43
WPA .....	43
The Right Hardware for the Task .....	44
Setting Up a Simple Wireless Network .....	44
An OpenBSD WPA Access Point .....	47
A FreeBSD WPA Access Point .....	48
The Access Point's PF Rule Set .....	49
Access Points with Three or More Interfaces .....	50
Handling IPSec, VPN Solutions .....	50
The Client Side .....	51
Guarding Your Wireless Network with <i>authpf</i> .....	54
A Basic Authenticating Gateway .....	55
Wide Open but Actually Shut .....	57

## 5 BIGGER OR TRICKIER NETWORKS 59

A Web Server and Mail Server on the Inside—Routable Addresses .....	60
A Degree of Separation: Introducing the DMZ .....	63
Sharing the Load: Redirecting to a Pool of Addresses .....	65
Getting Load Balancing Right with <i>relayd</i> .....	66
A Web Server and Mail Server on the Inside—the NAT Version .....	71
DMZ with NAT .....	73
Redirection for Load Balancing .....	73
Back to the Single NATed Network .....	74
Filtering on Interface Groups .....	76
The Power of Tags .....	77
The Bridging Firewall .....	78
Basic Bridge Setup on OpenBSD .....	79
Basic Bridge Setup on FreeBSD .....	80
Basic Bridge Setup on NetBSD .....	81
The Bridge Rule Set .....	82
Handling Nonroutable Addresses from Elsewhere .....	83

## **6 TURNING THE TABLES FOR PROACTIVE DEFENSE 85**

Turning Away the Brutes .....	86
SSH Brute-Force Attacks .....	86
Setting Up an Adaptive Firewall .....	86
Tidying Your Tables with pfctl .....	89
Giving Spammers a Hard Time with spamd .....	89
Network-Level Behavior Analysis and Blacklisting .....	90
Greylisting: My Admin Told Me Not to Talk to Strangers .....	93
Tracking Your Real Mail Connections: spamlogd .....	98
Greytrapping .....	98
Managing Lists with spamdb .....	100
Detecting Out-of-Order MX Use .....	102
Handling Sites That Do Not Play Well with Greylisting .....	102
Spam-Fighting Tips .....	104

## **7 QUEUES, SHAPING, AND REDUNDANCY 105**

Directing Traffic with ALTQ .....	105
Basic ALTQ Concepts .....	106
Queue Schedulers, aka Queue Disciplines .....	106
Setting Up ALTQ .....	107
Setting Up Queues .....	108
Priority-Based Queues .....	109
Class-Based Bandwidth Allocation for Small Networks .....	112
A Basic HFSC Traffic Shaper .....	113
Queueing for Servers in a DMZ .....	115
Using ALTQ to Handle Unwanted Traffic .....	117
Redundancy and Failover: CARP and pfsync .....	119
The Project Specification: A Redundant Pair of Gateways .....	119
Setting Up CARP .....	121
Keeping States Synchronized: Adding pfsync .....	125
Putting Together a Rule Set .....	126
CARP for Load Balancing .....	128

## **8 LOGGING, MONITORING, AND STATISTICS 131**

PF Logs: The Basics .....	132
Logging All Packets: log (all) .....	134
Logging to Several pflog Interfaces .....	135
Logging to Syslog, Local or Remote .....	135
Tracking Statistics for Each Rule with Labels .....	137
Additional Tools for PF Logs and Statistics .....	139
Keeping an Eye on Things with systat .....	139
Keeping an Eye on Things with pftop .....	141
Graphing Your Traffic with pfstat .....	141
Collecting NetFlow Data with pflow(4) .....	143
Collecting NetFlow Data with pfflowd .....	149
SNMP Tools and PF-Related SNMP MIBs .....	150
Log Data as the Basis for Effective Debugging .....	150

**GETTING YOUR SETUP JUST RIGHT****151**

Things You Can Tweak and What You Probably Should Leave Alone .....	151
Block Policy .....	152
Skip Interfaces .....	152
State Policy .....	153
State Defaults .....	153
Timeouts .....	154
Limits .....	155
Debug .....	156
Rule Set Optimization .....	157
Optimization .....	158
Fragment Reassembly .....	158
Cleaning Up Your Traffic .....	158
Packet Normalization with <code>scrub</code> .....	158
Protecting Against Spoofing with <code>antispoof</code> .....	159
Testing Your Setup .....	160
Debugging Your Rule Set .....	162
Know Your Network and Stay in Control .....	165

**RESOURCES****167**

General Networking and BSD Resources on the Internet .....	167
Sample Configurations and Related Musings .....	169
PF on Other BSD Systems .....	170
BSD and Networking Books .....	170
Wireless Networking Resources .....	171
<code>spamd</code> and Greylisting-Related Resources .....	171
Book-Related Web Resources .....	172
Buy OpenBSD CDs and Donate! .....	172

**A NOTE ON HARDWARE SUPPORT****173**

Getting the Right Hardware .....	174
Issues Facing Hardware Support Developers .....	175
How to Help the Hardware Support Efforts .....	175

**INDEX****177**

## **FOREWORD**

**from the first edition**

OpenBSD’s PF packet filter has enjoyed a lot of success and attention since it was first released in OpenBSD 3.0 in late 2001. While you’ll find out more about PF’s history in this book, in a nutshell, PF happened

because it was needed by the developers and users of OpenBSD. Since the original release, PF has evolved greatly and has become the most powerful free tool available for firewalling, load balancing, and traffic managing. When PF is combined with CARP and pfsync, PF lets system administrators not only protect their services from attack, but it makes those services more reliable by allowing for redundancy, and it makes them faster by scaling them using pools of servers managed through PF and relayd.

While I have been involved with PF’s development, I am first and foremost a large-scale user of PF. I use PF for security, to manage threats both internal and external, and to help me run large pieces of critical infrastructure in a redundant and scalable manner. This saves my employer (the University of Alberta, where I wear the head sysadmin hat by day) money, both in terms of downtime and in terms of hardware and software. You can use PF to do the same.

With these features comes the necessary evil of complexity. For someone well versed in TCP/IP and OpenBSD, PF's system documentation is quite extensive and usable all on its own. But in spite of extensive examples in the system documentation, it is never quite possible to put all the things you can do with PF and its related set of tools front and center without making the system documentation so large that it ceases to be useful for those experienced people who need to use it as a reference.

This book bridges the gap. If you are a relative newcomer, it can get you up to speed on OpenBSD and PF. If you are a more experienced user, this book can show you some examples of the more complex applications that help people with problems beyond the scope of the typical. For several years, Peter N.M. Hansteen has been an excellent resource for people learning how to apply PF in more than just the "How do I make a firewall?" sense, and this book extends his tradition of sharing that knowledge with others. Firewalls are now ubiquitous enough that most people have one, or several. But this book is not simply about building a firewall, it is about learning techniques for manipulating your network traffic and understanding those techniques enough to make your life as a system and network administrator a lot easier. A simple firewall is easy to build or buy off the shelf, but a firewall you can live with and manage yourself is somewhat more complex. This book goes a long way toward flattening out the learning curve and getting you thinking not only about how to build a firewall, but how PF works and where its strengths can help you. This book is an investment to save you time. It will get you up and running the right way—faster, with fewer false starts and less time experimenting.

**Bob Beck**  
**Director, The OpenBSD Foundation**  
*<http://www.openbsdfoundation.org>*  
**Edmonton, Alberta, Canada**

## **ACKNOWLEDGMENTS**

This manuscript started out as a user group lecture, first presented at the January 27, 2005 meeting of the Bergen [BSD and] Linux User Group (BLUG). After I had translated the manuscript into English and expanded it slightly, Greg Lehey suggested that I should stretch it a little further and present it as a half day tutorial for the AUUG 2005 conference. After a series of tutorial revisions, I finally started working on what was to become the book version in early 2007.

The next two paragraphs are salvaged from the tutorial manuscript and still apply to this book:

This manuscript is a slightly further developed version of a manuscript prepared for a lecture which was announced as (translated from Norwegian):

“This lecture is about firewalls and related functions, with examples from real life with the OpenBSD project’s PF (Packet Filter). PF offers firewalling, NAT, traffic control, and bandwidth management in a single, flexible, and sysadmin-friendly system. Peter hopes that the lecture will give you some ideas about how to

control your network traffic the way you want—keeping some things outside your network, directing traffic to specified hosts or services, and of course, giving spammers a hard time.”

Some portions of content from the tutorial (and certainly all the really useful topics) made it into this book in some form. During the process of turning it into a useful book, a number of people have offered insights and suggestions.

People who have offered significant and useful input regarding early versions of this manuscript include Eystein Roll Aarseth, David Snyder, Peter Postma, Henrik Kramshøj, Vegard Engen, Greg Lehey, Ian Darwin, Daniel Hartmeier, Mark Uemura, Hallvor Engen, and probably a few who will remain lost in my mail archive until I can grep them out of there.

I would like to thank the following organizations for their kind support: the NUUG Foundation for a travel grant, which partly financed my AUUG 2005 appearance; the AUUG, UKUUG, SANE, BSDCan, and AsiaBSDCon organizations for inviting me to their conferences; and the FreeBSD Foundation for sponsoring my trips to BSDCan 2006 and EuroBSDCon 2006.

Much like the first, the second edition was written mainly at night and on weekends, as well as during other stolen moments at odd hours. I would like to thank my former colleagues at FreeCode for easing the load for a while by allowing me some chunks of time to work on the second edition in between other projects during the early months of 2010. I would also like to thank several customers, who have asked that their names not be published, for their interesting and challenging projects, which inspired some of the configurations offered here. You know who you are.

Finally, during the process of turning the manuscript into a book, several people did amazing things that helped this book become a lot better. I am indebted to Bill Pollock and Adam Wright for excellent developmental editing; I would like to thank Henning Brauer for excellent technical review; heartfelt thanks go to Eystein Roll Aarseth, Jakob Breivik Grimstveit, Hallvor Engen, Christer Solskogen, Ian Darwin, Jeff Martin, and Lars Noodén for valuable input on various parts of the manuscript; and, finally, warm thanks to Megan Dunchak and Linda Recktenwald for their efforts in getting the first edition of the book into its final shape and to Serena Yang for guiding the second edition to completion. Special thanks are due to Dru Lavigne for making the introductions which led to this book getting written in the first place, instead of just hanging around as an online tutorial and occasional conference material.

Last but not least, I would like to thank my dear wife, Birthe, and my daughter, Nora, for all their love and support, before and during the book writing process as well as throughout the rather intense work periods that yielded the second edition. This would not have been possible without you.

# INTRODUCTION



This is a book about building the network you need. We'll dip into the topics of firewalls and related functions, starting from *a little* theory. You'll see plenty of examples of filtering and other ways to direct network traffic. I'll assume that you have a basic to intermediate command of TCP/IP networking concepts and Unix administration.

All the information in this book comes with a fair warning: As in any number of other endeavors, the solutions we discuss can be done in *more than one way*. You should also be aware that the software world could have changed slightly or quite a bit since the book was printed.

The information in the book is as up to date and correct as possible at the time of writing, and refers to OpenBSD version 4.8, FreeBSD 8.1, and NetBSD 5.0, with any patches available in late August 2010.

## This Is Not a HOWTO

The book is a direct descendant of a moderately popular PF tutorial. The tutorial is also the source of the following admonition, and you may be exposed to this live if you attend one of my tutorial sessions:

This document is not intended as a precooked recipe for cutting and pasting.

Just to hammer this in, please repeat after me:

---

The Pledge of the Network Admin

This is my network.

It is mine,  
or technically, my employer's.  
It is my responsibility,  
and I care for it with all my heart.

There are many other networks a lot like mine,  
but none are just like it.

I solemnly swear

that I will not mindlessly paste from HOWTOs.

---

The point is that while the configurations I show you do work (I have tested them, and they are in some way related to what has been put into production), they may be overly simplistic, since many were designed to demonstrate a specific point of configuration. They are almost certain to be at least a little off, and they possibly could be quite wrong for your network.

Please keep in mind that this book is intended to show you a few useful techniques and inspire you to achieve good things.

Please strive to understand your network and what you need to do to make it better.

Please do not paste blindly from this document or any other.

## What This Book Covers

The book is intended to be a stand-alone document to enable you to work on your machines with only short forays into man pages and occasional reference to the online and printed resources listed in Appendix A.

Your system probably comes with a prewritten *pf.conf* file containing some commented-out suggestions for useful configurations, as well as a few examples in the documentation directories such as */usr/share/pf/*. These examples are useful as a reference, but we won't use them directly in this book. Instead, you'll learn how to construct a *pf.conf* from scratch, step by step.

Here is a brief rundown of what you will find in this book:

- Chapter 1, “Building the Network You Need,” walks through basic networking concepts, gives a short overview of PF’s history, and provides some pointers on how to adjust to the BSD way if you are new to this family of operating systems. Read this chapter first if you want to get your general bearings for working with BSD systems.
- Chapter 2, “PF Configuration Basics,” shows you how to enable PF on your system and covers a very basic rule set for a single machine. This chapter is a fairly crucial one, since all the later configurations are based on the one we build in this chapter.
- Chapter 3, “Into the Real World,” builds on the single-machine configuration in Chapter 2 and leads you through the basics of setting up a gateway that serves as a point of contact between separate networks. By the end of Chapter 3, you’ll have built a configuration that is fairly typical for a home or small office network, with some tricks up your sleeve to make network management easier. You’ll also get an early taste of how to handle services with odd requirements such as FTP, as well as some tips on how to make your network troubleshooting-friendly by catering to some of the frequently less understood Internet protocols and services.
- Chapter 4, “Wireless Networks Made Easy,” walks you through adding wireless networking to your setup. The wireless environment presents some security challenges, and by the end of this chapter, you may find yourself with a wireless network with access control and authentication via authpf. Some of the information is likely to be useful in wired environments, too.
- Chapter 5, “Bigger or Trickier Networks,” tackles the situation where you introduce servers and services that need to be accessible from outside your own network. By the end of this chapter, you may have a network with one or several separate subnets and DMZs, and you will have tried your hand at a couple of different load-balancing schemes via redirections and relayd in order to improve service quality for your users.
- Chapter 6, “Turning the Tables for Proactive Defense,” shows you some of the tools in the PF tool chest for dealing with attempts at undesirable activity, and how to use them productively. Here, we deal with brute-force password-guessing attempts and other network flooding, as well as the ever-favorite antispam tool spamd, the OpenBSD spam deferral daemon. This chapter should make your network a more pleasant one for legitimate users and not so welcoming to those with less than good intentions.
- Chapter 7, “Queues, Shaping, and Redundancy,” introduces traffic shaping via the ALTQ queueing engine. We then move on to creating redundant configurations, with CARP configurations for both failover and load balancing. This chapter should leave you with better resource utilization through traffic shaping adapted to your network needs, as well as better availability with a redundant, CARP-based configuration.

- Chapter 8, “Logging, Monitoring, and Statistics,” explains PF logs. You’ll learn how to extract and process log and statistics data from your PF configuration with tools in the base system as well as optional packages. This is where you will be exposed to NetFlow and SNMP-based tools.
- Chapter 9, “Getting Your Setup Just Right,” walks through various options that will help you tune your setup. It ties together the knowledge you have gained from the previous chapters with a rule set debugging tutorial.
- Appendix A, “Resources,” is an annotated list of print and online literature and other resources you may find useful as you expand your knowledge of PF and networking topics.
- Appendix B, “A Note on Hardware Support,” gives an overview of some of the issues involved in creating a first-rate tool as free software.

If you’re confident in your skills, you can jump to the chapter or section that interests you the most. However, each successive chapter builds on work done in the earlier chapters, so it may be useful to read through the chapters in sequence. The main perspective in the book is the world as seen from the command line in OpenBSD 4.8, with notes on other systems where there are significant differences.

# 1

## BUILDING THE NETWORK YOU NEED



PF, the OpenBSD *Packet Filter subsystem*, is one of the finest tools available for taking control of your network. Before diving into the specifics of how to make your network the fine-tuned machinery of your dreams, please read this chapter. It introduces basic networking terminology and concepts, provides some PF history, and gives you an overview of what you can expect to find in this book.

### **Your Network: High Performance, Low Maintenance, and Secure**

If this heading accurately describes your network, you're most likely reading this for pure entertainment, and I hope you will enjoy the rest of the book. If, on the other hand, you're still learning how to build networks or you're

not quite confident of your skills yet, a short recap of basic network security concepts can be useful.

Information technology (IT) security is a large, complex and sometimes confusing subject. Even if we limit ourselves to thinking only in terms of network security, there is a perception that we haven't really narrowed down the field much or eliminated enough of the inherently confusing terminology. Matters became significantly worse some years ago when personal computers started joining the networked world, equipped with system software and applications that were clearly not designed for a networked environment.

The result was rather predictable. Even before the small computers became networked, they had become home to malicious software such as *viruses* (semiautonomous software that is able to "infect" other files in order to deliver its payload and make further copies of itself) and *trojans* (originally *trojan horses*, software or documents with code embedded that if activated would cause the victim's computer to perform actions that the user did not intend). When the small computers became networked, they were introduced to yet another kind of malicious software called a *worm*, a class of software that uses the network to propagate its payload.<sup>1</sup> Along the way, the networked versions of various kinds of frauds made it onto the network security horizon as well, and today a significant part of computer security activity (possibly the largest segment of the industry) centers on threat management, with emphasis on fighting and cataloging malicious software, or *malware*.

The futility of enumerating badness has been argued convincingly elsewhere (see Appendix A for references, such as Marcus Ranum's excellent essay "The Six Dumbest Ideas in Computer Security"). The OpenBSD approach is to design and code properly in the first place. Then if you later discover mistakes, and the bugs turn out to be exploitable, fix those bugs everywhere similar code turns up in the tree, even if it could mean a radical overhaul of the design and, at worst, a loss of backward compatibility.<sup>2</sup>

In PF, and by extension in this book, the focus is narrower, concentrated on network traffic at the network level. The introduction of `divert(4)` sockets in OpenBSD 4.7 made it incrementally easier to set up a system where PF contributes to *deep packet inspection*, much like some fiercely marketed products. However, no widely used free software yet uses the interface, and we will instead focus on some techniques based on pure network-level behavior (most evident in the example configurations in Chapter 6) that will help ease the load on the content-inspecting products if you have them in place. As you will see in the following chapters, the network level offers a lot of fun and excitement, in addition to the blocking or passing packets.

---

1. The famous worms before the Windows era were the IBM Christmas Tree EXEC worm (1987) and the first Internet worm, the Morris worm (1988), both within easy reach of your favorite search engine. The Windows era of networked worms is considered to have started with the ILOVEYOU worm in May 2000.

2. Several presentations on OpenBSD's approach to security can be found via <http://www.openbsd.org/papers/>. Some of my favorites are Theo de Raadt's "Exploit Mitigation Techniques," Damien Miller's "Security Measures in OpenSSH," and "Puffy at Work—Getting Code Right and Secure, the OpenBSD Way," by Henning Brauer and Sven Dehmlow.

## Where the Packet Filter Fits In

The packet filter's main function is, as the name suggests, to filter network packets by matching the properties of individual packets and the network connections built from those packets against the filtering criteria defined in its configuration files. The packet filter is responsible for deciding what to do with those packets. That could mean passing them through or rejecting them, or triggering events that other parts of the operating system or external applications are set up to handle.

PF lets you write custom filtering criteria to control network traffic based on essentially any packet or connection property, including address family, source and destination address, interface, protocol, port, and direction. Based on these criteria, the packet filter performs the action you specify. One of the simplest and most common actions is to block traffic.

A packet filter can keep unwanted traffic out of your network. It can also help contain network traffic inside your own network. Both those functions are important to the *firewall* concept, but blocking is far from the only useful or interesting feature of a functional packet filter. As you will see in this book, you can use filtering criteria to direct certain kinds of network traffic to specific hosts, assign classes of traffic to queues, perform traffic shaping, and even hand off selected kinds of traffic to other software for special treatment.

All this processing happens at the network level, based on packet and connection properties. PF is part of the network stack, firmly embedded in the operating system kernel. While there have been examples of packet filtering implemented in user space, in most operating systems, the filtering functions are performed in the kernel because it's faster to do so.

## The Rise of PF

If you have a taste for history, you probably already know that OpenBSD and the other BSDs<sup>3</sup> are direct descendants of the BSD system (sometimes referred to as *BSD Unix*), the operating system that contained the original reference implementation of the TCP/IP Internet protocols in the early 1980s.

As the research project behind BSD development started winding down in the early 1990s, the code was liberated for further development by small groups of enthusiasts around the world. Some of these enthusiasts were responsible for keeping vital parts of the emerging Internet's infrastructure running reliably, and BSD development continued along parallel lines in

3. If *BSD* does not sound familiar, here is a short explanation. The acronym expands to *Berkeley Software Distribution* and originally referred to a collection of useful software developed for the Unix operating system by staff and students at the University of California, Berkeley. Over time, the collection expanded into a complete operating system, which in turn became the forerunner of a family of systems, including OpenBSD, FreeBSD, NetBSD, DragonFly BSD, and, by some definitions, even Apple's Mac OS X. For a very readable explanation of what BSD is, see Greg Lehey's "Explaining BSD" at <http://www.freebsd.org/doc/en/articles/explaining-bsd> (and, of course, the projects' websites).

several groups. The OpenBSD group became known as the most security-oriented of the BSDs. For its packet filtering needs, it used a subsystem called IPFilter, written by Darren Reed.

It shocked the OpenBSD community when Reed announced in early 2001 that IPFilter, which at that point was intimately integrated with OpenBSD, was not covered under the BSD license. Instead, it used almost a word-for-word copy of the license, omitting only the right to make changes to the code and distribute the result. The problem was that the OpenBSD version of IPFilter contained several changes and customizations, which, as it turned out, were not allowed under the license. As a result, IPFilter was removed from the OpenBSD source tree on May 29, 2001, and for a few weeks, the development version of OpenBSD (-current) did not include any firewalling software.

Fortunately, at this time, in Switzerland, Daniel Hartmeier had been performing some limited experiments involving kernel hacking in the networking code. He began by hooking a small function of his own into the networking stack and then making packets pass through it. Then he began thinking about filtering. When the license crisis happened, PF was already well under development. The first commit of the PF code was on Sunday, June 24, 2001, at 19:48:58 UTC. A few months of intense activity followed, and the resulting version of PF was launched as a default part of the OpenBSD 3.0 base system in December of 2001.<sup>4</sup> This version contained a rather complete implementation of packet filtering, including network address translation, with a configuration language that was similar enough to IPFilter's that migrating to the new OpenBSD version did not pose major problems.<sup>5</sup>

PF proved to be well-developed software. In 2002, Hartmeier presented a USENIX paper with performance tests showing that the OpenBSD 3.1 PF performed equally well or better under stress than either IPFilter on OpenBSD 3.1 or iptables on Linux. In addition, tests run on the original PF from OpenBSD 3.0 showed mainly that the code had gained in efficiency from version 3.0 to version 3.1.<sup>6</sup>

The OpenBSD PF code, with a fresh packet-filtering engine written by experienced and security-oriented developers, naturally generated interest in the sister BSDs as well. The FreeBSD project gradually adopted PF, first as

---

4. The IPFilter copyright episode spurred the OpenBSD team to perform a license audit of the entire source tree and ports in order to avoid similar situations in the future. Several potential problems were resolved over the months that followed, resulting in the removal of a number of potential license pitfalls for everyone involved in free software development. Theo de Raadt summed up the effort in a message to the *openbsd-misc* mailing list on February 20, 2003. The initial drama of the license crisis had blown over, and the net gain was a new packet filtering system under a free license, with the best code quality available, as well as better free licenses for a large body of code in OpenBSD itself and in other widely used free software.

5. Compatibility with IPFilter configurations was an early design goal for the PF developers, but it stopped being a priority once it could be safely assumed that all OpenBSD users had moved to PF (around the time OpenBSD 3.2 was released, if not earlier). You should not assume that an existing IPFilter configuration will work without changes with any version of PF. With the syntax changes introduced in OpenBSD 4.7, even upgrades from earlier PF versions will involve some conversion work.

6. The article that provides the details of these tests is available from Daniel Hartmeier's website. See <http://www.benzedrine.cx/pf-paper.html>.

a package, and then from version 5.3 on, in the base system as one of three packet filtering systems. PF has also been included in NetBSD and DragonFly BSD.<sup>7</sup>

### NEWER PF RELEASES PERFORM BETTER

Like the rest of the computing world, OpenBSD and PF have been affected by rapid changes in hardware and network conditions. I have not seen comparable tests to the ones in Daniel Hartmeier's USENIX paper performed recently, but PF users have found that the PF filtering overhead is rather modest.

As an example (mainly to illustrate that even unexciting hardware configurations can be useful), the machine that gateways between one small office network in my care and the world is a Pentium III 450MHz with 384MB of RAM. When I've remembered to check, I've never seen the machine at less than 96 percent idle according to top.

It is also worth noting that the current PF developers, mainly Henning Brauer and Ryan McBride, with contributions from several others, have introduced a number of optimizations to OpenBSD's PF code during recent releases, making each release from 4.4 through 4.8 perform noticeably better than its predecessors.

This book focuses on the PF version available in OpenBSD 4.8. I will note significant differences between that version and the ones integrated in other systems as appropriate.

If you're ready to dive into PF configuration, you can jump to Chapter 2 to get started. If you want to spend a little more time getting your bearings in unfamiliar BSD territory, continue reading this chapter.

## If You Came from Elsewhere

If you are reading this because you are considering moving your setup to PF from some other system, this section is for you.

If you want to use PF, you need to install and run a BSD system such as OpenBSD, FreeBSD, NetBSD, or DragonFly BSD. These are all fine operating systems, but my personal favorite is OpenBSD, mainly because that is the operating system where essentially all PF development happens, and I find the developers' and the system's no-nonsense approach refreshing.

Occasionally, minor changes and bug fixes trickle back to the main PF code base from the PF implementations on other systems, but the newest, most up-to-date PF code is always to be found on OpenBSD. Some of the features described in this book are available only in the most recent versions of OpenBSD. The other BSDs tend to port the latest released PF version from OpenBSD to their code bases in time for their next release, but synchronized updates are far from guaranteed, and the lag is sometimes considerable.

<sup>7</sup>. At one point even a personal firewall product, Core Force, claimed to be based on PF. By early 2010, Core Security, the company that developed Core Force (<http://force.coresecurity.com/>), seemed to have shifted focus to other security areas such as penetration testing, but the product was still available for download.

If you are planning to run PF on FreeBSD, NetBSD, DragonFly BSD, or another system, you should check your system’s release notes and other documentation for information about which version of PF is included.

## Pointers for Linux Users

The differences and similarities between Linux and BSD are potentially a large topic if you probe deeply, but if you have a reasonable command of the basics, it should not take too long for you to feel right at home with the BSD way of doing things. In the rest of this book, I will assume that you can find your way around the basics of BSD network configuration. So, if you are more familiar with configuring Linux or other systems than you are with BSD, it is worth noting a few points about BSD configuration.

- Linux and BSD use different conventions for naming network interfaces. The Linux convention is to label all the network interfaces on a given machine in the sequence eth0, eth1, and so on (although with some Linux versions and driver combinations, you also see wlan0, wlan1, and so on for wireless interfaces).

On the BSDs, interfaces are assigned names that equal the driver name plus a sequence number. For example, older 3Com cards using the ep driver appear as ep0, ep1, and so on; Intel Gigabit cards are likely to end up as em0, em1, and so on. Some SMC cards are listed as sn0, sn1, and so on. This system is quite logical, and makes it easier to find the documentation for the specifics of that interface. If your kernel reports (at boot time or in ifconfig output) that you have an interface called em0, you need only type `man em` at a shell command-line prompt to find out what speeds it supports, whether there are any eccentricities to be aware of, whether any firmware download is needed, and so on.
- You should be aware that in BSDs, the configuration is `/etc/rc.conf`-centric. In general, the BSDs are organized to read the configuration from the file `/etc/rc.conf`, which is read by the `/etc/rc` script at startup. OpenBSD recommends using `/etc/rc.conf.local` for local customizations, since `rc.conf` contains the default values. FreeBSD uses `/etc/defaults/rc.conf` to store the default settings, making `/etc/rc.conf` the correct place to make changes. In addition, OpenBSD uses per-interface configuration files called `hostname.<if>`, where `<if>` is replaced with the interface name.
- For the purpose of learning PF, you will need to concentrate on an `/etc/pf.conf` file, which will be largely your own creation.

If you need a broader and more thorough introduction to your BSD of choice, look up the operating system’s documentation, including FAQs and guides, at the project’s website. You can also find some suggestions for further reading in Appendix A.

## **Frequently Answered Questions About PF**

This section is based on questions I've been asked via email or at meetings and conferences, as well as some that have popped up in mailing lists and other discussion forums. Some of the more common questions are covered here, in a FAQ-style<sup>8</sup> format.

### **Can I run PF on my Linux machine?**

In a word, no. Over the years, announcements have appeared on the PF mailing list from someone claiming to have started a Linux port of PF, but at the time of this writing, no one has yet claimed to have completed the task. The main reason for this is probably that PF is developed primarily as a deeply integrated part of the OpenBSD networking stack. Even after more than a decade of parallel development, the OpenBSD code still shares enough fundamentals with the other BSDs to make porting possible, but porting PF to a non-BSD system would require rewriting large chunks of PF itself, as well as whatever integration is needed at the target side.

For some basic orientation tips for Linux users to find their way in BSD network configurations, see “Pointers for Linux Users” on page 6.

### **Can you recommend a GUI tool for managing my PF rule set?**

This book is mainly oriented toward users who edit their rule sets in their favorite text editor. The sample rule sets in this book are simple enough that you probably would not get a noticeable benefit from any of the visualization options the various GUI tools are known to offer.

A rather common claim is that the PF configuration files are generally readable enough that a graphic visualization tool is not really necessary. There are, however, several GUI tools available that can edit and/or generate PF configurations, including a complete, customized build of FreeBSD called pfSense (<http://www.pfsense.org/>), which includes a sophisticated GUI rule editor.

I recommend that you work through the parts of this book that apply to your situation, and then decide if you need to use a GUI tool to feel comfortable running and maintaining the systems you build.

### **Is there a tool I can use to convert my OtherProduct® setup to a PF configuration?**

The best strategy when converting network setups, including firewall setups, from one product to another is to go back to the specifications or policies for your network or firewall configuration, and then implement the policies using the new tool.

Other products will inevitably have a slightly different feature set, and the existing configuration you created for OtherProduct® is likely to mirror slightly different approaches to specific problems, which do not map easily, or at all, to features in PF and related tools.

---

<sup>8</sup>. The three-letter abbreviation *FAQ* expands to either *frequently asked questions* or *frequently answered questions*—both equally valid.

Having a documented policy, and taking care to update it as your needs change, will make your life easier. This documentation should contain a complete prose specification of what your setup is meant to achieve. (You might start out by putting comments in your configuration file to explain the purpose of your rules.) This makes it possible to verify whether the configuration you are running actually implements the design goals. In some corporate settings, there may even be a formal requirement for a written policy.

The impulse to look for a way to automate your conversion is quite understandable and perhaps expected in a system administrator. I urge you to resist the impulse and to perform your conversion after reevaluating your business and technical needs and (preferably) after creating or updating a formal specification or policy in the process.

Some of the GUI tools that serve as administration front ends claim the ability to output configuration files for several firewall products, and could conceivably be used as conversion tools. However, this has the effect of inserting another layer of abstraction between you and your rule set, and puts you at the mercy of the tool author's understanding of how PF rule sets work. I recommend working through at least the relevant parts of this book before spending serious time on considering an automated conversion.

### **Why did the PF rules syntax change all of a sudden?**

The world changed, and PF changed with it. More specifically, the OpenBSD developers have a very active and pragmatically critical relationship to their code, and like all parts of OpenBSD, the PF code is under constant review.

The lessons learned over almost a decade of PF development and use led to internal changes in the code that eventually made it clear to the developers that changing the syntax slightly would make sense. The result for you, the user, is that PF is now even easier to use and performs better than the earlier versions. If you are upgrading your system to OpenBSD 4.7 or newer, you are in for a real treat.

### **Where can I find out more?**

There are several good sources of information about PF and the systems on which it runs. You have already found one in this book. You can find references to a number of printed and online resources in Appendix A.

If you have a BSD system with PF installed, consult the online manual pages (*man pages*) for information about your exact release of the software. Unless otherwise indicated, the information in this book refers to the world as it looks from the command line on an OpenBSD 4.8 system.

## A Little Encouragement: A PF Haiku

If you are not quite convinced yet (or even if you are reading on anyway), a little encouragement may be in order. Over the years, a good many people have said and written their bit about PF—sometimes odd, sometimes wonderful, and sometimes just downright strange.

The poem quoted here is a good indication of the level of feeling PF sometimes inspires in its users. This poem appeared on the PF mailing list, in a thread that started with a message with the subject “Things pf can’t do?” in May 2004. The message was written by someone who did not have a lot of firewall experience, and who consequently found it hard to get the desired setup.

This, of course, led to some discussion, with several participants saying that if PF was hard on a newbie, the alternatives were certainly not a bit better. The thread ended in the following haiku of praise from Jason Dixon, dated May 20, 2004.

---

Compared to working with iptables, PF is like this haiku:

A breath of fresh air,  
floating on white rose petals,  
eating strawberries.

Now I'm getting carried away:

Hartmeier codes now,  
Henning knows not why it fails,  
fails only for noob.

Tables load my lists,  
tarpit for the asshole spammer,  
death to his mail store.

CARP due to Cisco,  
redundant blessed packets,  
licensed free for me.

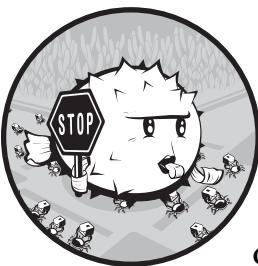
---

Some of the concepts Dixon mentions here may sound a bit unfamiliar, but if you read on, it will all make sense soon.



# 2

## PF CONFIGURATION BASICS



In this chapter, we will create a very simple setup with PF. We'll begin with the simplest configuration possible: a single machine configured to communicate with a single network. That network could very well be the Internet.

Your two main tools for configuring PF are your favorite text editor and the `pfctl` command-line administration tool. PF configurations, usually stored in `/etc/pf.conf`, are called *rule sets*, because each line in the configuration file is a *rule* that helps determine what the packet-filtering subsystem should do with the network traffic it sees. In ordinary, day-to-day administration, you will edit your configuration in the `/etc/pf.conf` file, and then load your changes using `pfctl`. There are web interfaces for PF administration tasks, but they are not part of the base system. The PF developers are not hostile toward these options, but they have yet to see a graphical interface for configuring PF that is clearly preferable to editing `pf.conf` and using `pfctl` commands.

## The First Step: Enabling PF

Before you can get started on the fun parts of shaping your network with PF and related tools, you need to make sure that PF is available and enabled. The details depend on your specific operating system: OpenBSD, FreeBSD, or NetBSD. Check your setup by following the instructions for your operating system, and then move on to “A Simple PF Rule Set: A Single, Stand-Alone Machine” on page 16.

The `pfctl` command is a program that requires higher privilege than the default for ordinary users. In the rest of this book, you will see commands that require extra privilege prefixed with `sudo`. If you have not started using `sudo` yet, you should. `sudo` is in the base system on OpenBSD. On FreeBSD and NetBSD, it is within easy reach via the `ports` system or `pkgsrc` system, respectively, as *security/sudo*.

Here are a couple general notes regarding using `pfctl`:

- The command to disable PF is `pfctl -d`. Once you have entered that command, all PF-based filtering that may have been in place will be disabled, and all traffic will be allowed to pass.
- For convenience, `pfctl` can handle several operations on a single command line. To enable PF and load the rule set in a single command, enter this:

---

```
$ sudo pfctl -ef /etc/pf.conf
```

---

### Setting Up PF on OpenBSD

In OpenBSD 4.6 and later, you do not need to enable PF, since it is enabled by default with a minimal configuration in place.<sup>1</sup> If you were watching the system console closely while the system was starting up, you may have noticed the `pf enabled` message appear soon after the kernel messages completed.

If you did not see the `pf enabled` message on the console at startup, you have several options to check that PF is indeed enabled. One simple way to check is to enter the command you would otherwise use to enable PF from the command line, like this:

---

```
$ sudo pfctl -e
```

---

If PF is already enabled, the system responds with this message:

---

```
pfctl: pf already enabled
```

---

If PF is not enabled, the `pfctl -e` command will enable PF and display this:

---

```
pf enabled
```

---

1. If you are setting up your first PF configuration on an OpenBSD version earlier than this, the best advice is to upgrade to the most recent stable version. If for some reason you must stay with the older version, it might be useful to consult the first edition of this book, as well as the man pages and other documentation for the specific version you are using.

In versions prior to OpenBSD 4.6, PF was not enabled by default. You can override the default by editing your `/etc/rc.conf.local` file (or creating the file, if it does not exist). Although it is not necessary on recent OpenBSD versions, it does not hurt to add this line to your `/etc/rc.conf.local` file:

---

pf=YES	# enable PF
--------	-------------

---

If you take a look at the `/etc/pf.conf` file in a fresh OpenBSD installation, you get your first exposure to a working rule set.

The default OpenBSD `pf.conf` file starts off with a `set skip on lo` rule to make sure traffic on the loopback interface is not filtered in any way. The next active line is a simple `pass` default to let your network traffic pass by default. Finally, an explicit `block` rule blocks remote X11 traffic to your machine.

As you probably noticed, the default `pf.conf` file also contains a few comment lines starting with a hash mark (#). In those comments, you will find suggested rules that hint at useful configurations such as FTP proxying (see Chapter 3) and `spamd`, the OpenBSD spam-deferral daemon (see Chapter 6). These items are potentially useful in various real-world scenarios, but since they may not be relevant in all configurations, they are commented out in the file by default.

If you look for PF-related settings in your `/etc/rc.conf` file, you will find the setting `pf_rules=`. In principle, this lets you specify that your configuration is in a file other than the default `/etc/pf.conf`. However, changing this setting is probably not worth the trouble. Using the default setting lets you take advantage of a number of automatic housekeeping features, such as automatic nightly backup of your configuration to `/var/backups`.

On OpenBSD, the `/etc/rc` script has a built-in mechanism to help you out if you reboot with either no `pf.conf` file or one that contains an invalid rule set. Before enabling any network interfaces, the `rc` script loads a rule set that allows a few basic services: SSH from anywhere, basic name resolution, and NFS mounts. This allows you to log in and correct any errors in your rule set, load the corrected rule set, and then go on working from there.

## **Setting Up PF on FreeBSD**

Good code travels well, and FreeBSD users will tell you that good code from elsewhere tends to find its way into FreeBSD sooner or later. PF is no exception, and from FreeBSD 5.2.1 and the 4.x series onward, PF and related tools became part of FreeBSD.

If you read through the previous section on setting up PF on OpenBSD, you saw that on OpenBSD, PF is enabled by default. That is not the case on FreeBSD, where PF is one of three possible packet-filtering options. Here, you need to take explicit steps to enable PF, and compared to OpenBSD, it

seems that you need a little more magic in your */etc/rc.conf*. A look at your */etc/default/rc.conf* file shows that the FreeBSD default values for PF-related settings are as follows:

---

```
pf_enable="NO"          # Set to YES to enable packet filter (pf)
pf_rules="/etc/pf.conf"  # rules definition file for pf
pf_program="/sbin/pfctl" # where pfctl lives
pf_flags=""             # additional flags for pfctl
pflog_enable="NO"        # set to YES to enable packet filter logging
pflog_logfile="/var/log/pflog" # where pflogd should store the logfile
pflog_program="/sbin/pflogd" # where pflogd lives
pflog_flags=""           # additional flags for pflogd
pfsync_enable="NO"        # expose pf state to other hosts for syncing
pfsync_syncdev=""         # interface for pfsync to work through
pfsync_ifconfig=""        # additional options to ifconfig(8) for pfsync
```

---

Fortunately, you can safely ignore most of these—at least for now. The following are the only options that you need to add to your */etc/rc.conf* configuration:

---

```
pf_enable="YES"          # Enable PF (load module if required)
pflog_enable="YES"        # start pflogd(8)
```

---

There are some differences between the FreeBSD 4.x, 5.x, and later with respect to PF. Refer to the *FreeBSD Handbook*, specifically the PF section of the “Firewalls” chapter, available from <http://www.freebsd.org/>, to see which information applies in your case. The PF code in FreeBSD 7 and 8 is equivalent to the code in OpenBSD 4.1 with some bug fixes. The instructions in this book assume that you are running FreeBSD 7.0 or newer.

On FreeBSD, PF is compiled as a kernel-loadable module by default. If your FreeBSD setup runs with a GENERIC kernel, you should be able to start PF with the following:

---

```
$ sudo kldload pf
$ sudo pfctl -e
```

---

Assuming you have put the lines just mentioned in your */etc/rc.conf* and created an */etc/pf.conf* file, you could also use the PF *rc* script to run PF. The following enables PF:

---

```
$ sudo /etc/rc.d/pf start
```

---

And this disables the packet filter:

---

```
$ sudo /etc/rc.d/pf stop
```

---

**NOTE** On FreeBSD, the `/etc/rc.d/pf` script requires at least a line in `/etc/rc.conf` that reads `pf_enable="YES"` and a valid `/etc/pf.conf` file. If either of these requirements is not met, the script will exit with an error message. There is no `/etc/pf.conf` file in a default FreeBSD installation, so you'll need to create one before you reboot the system with PF enabled. For our purposes, creating an empty `/etc/pf.conf` with `touch` will do, but you could also work from a copy of the `/usr/share/examples/pf/pf.conf` file supplied by the system.

The supplied sample file `/usr/share/examples/pf/pf.conf` contains no active settings. It has only comment lines starting with a # character and commented-out rules, but it does give you a preview of what a working rule set will look like. For example, if you remove the # sign before the line that says `set skip on lo` to uncomment the line, and then save the file as your `/etc/pf.conf`, once you enable PF and load the rule set, your loopback interface will not be filtered. However, even if PF is enabled on your FreeBSD system, we haven't gotten around to writing an actual rule set, so PF is not doing much of anything and all packets will pass.

As of this writing (September 2010), the FreeBSD `rc` scripts do not set up a default rule set as a fallback if the configuration read from `/etc/pf.conf` fails to load. This means that enabling PF with no rule set or with `pf.conf` content that is syntactically invalid will leave the packet filter enabled with a default `pass all` rule set.

## Setting Up PF on NetBSD

On NetBSD 2.0, PF became available as a loadable kernel module that could be installed via packages (`security/pf.lkm`) or compiled into a static kernel configuration. In NetBSD 3.0 and later, PF is part of the base system. On NetBSD, PF is one of several possible packet-filtering systems, and you need to take explicit action to enable it.

Some details of PF configuration have changed between NetBSD releases. This book assumes you are using NetBSD 5.0 or later.<sup>2</sup>

To use the loadable PF module for NetBSD, add the following lines to your `/etc/rc.conf` to enable loadable kernel modules, PF, and the PF log interface, respectively.

---

```
lkm="YES" # do load kernel modules
pf=YES
pflogd=YES
```

---

To load the `pf` module manually and enable PF, enter this:

---

```
$ sudo modload /usr/lkm/pf.o
$ sudo pfctl -e
```

---

2. For instructions on using PF in earlier releases, see the documentation for your release and look up supporting literature listed in Appendix A of this book.

Alternatively, you can run the *rc.d* scripts to enable PF and logging, as follows:

---

```
$ sudo /etc/rc.d/pf start  
$ sudo /etc/rc.d/pflogd start
```

---

To load the module automatically at startup, add the following line to */etc/lkm.conf*:

---

```
/usr/lkm/pf.o - - - - BEFORENET
```

---

If your */usr* filesystem is on a separate partition, add this line to your */etc/rc.conf*:

---

```
critical_filesystems_local="${critical_filesystems_local} /usr"
```

---

If there are no errors at this point, you have enabled PF on your system, and you are ready to move on to creating a complete configuration.

The supplied */etc/pf.conf* file contains no active settings. It has only comment lines starting with a hash mark (#) and commented-out rules, but it does give you a preview of what a working rule set will look like. For example, if you remove the hash mark before the line that says `set skip on lo` to uncomment it, and then save the file, once you enable PF and load the rule set, your loopback interface will not be filtered. However, even if PF is enabled on your NetBSD system, we haven't gotten around to writing an actual rule set, so PF is not doing much of anything but passing packets.

NetBSD implements a default or fallback rule set via the file */etc/defaults/pf.boot.conf*. This rule set is intended only to let your system complete its boot process in case the */etc/pf.conf* file does not exist or contains an invalid rule set. You can override the default rules by putting your own customizations in */etc/pf.boot.conf*.

## A Simple PF Rule Set: A Single, Stand-Alone Machine

Mainly to have a common, minimal baseline, we will start building rule sets from the simplest possible configuration.

### A Minimal Rule Set

The simplest possible PF setup is on a single machine that will not run any services and talks to only one network (which may be the Internet).

We'll begin with an */etc/pf.conf* file that looks like this:

---

```
block in all  
pass out all keep state
```

---

This rule set denies all incoming traffic, allows traffic we send, and retains state information on our connections. PF reads rules from top to bottom; the *last* rule in a rule set that matches a packet or connection is the one that is applied.

Here, any connection coming into our system from anywhere else will match the `block in all` rule. Even with this tentative result, the rule evaluation will continue to the next rule (`pass out all keep state`), but the traffic will not even match the first criterion (the `out` direction) in this rule. With no more rules to evaluate, the status will not change, and the traffic will be blocked. In a similar manner, any connection initiated from the machine with this rule set will not match the first rule, but will match the second (once again, the wrong direction); matching a `pass` rule, it is allowed to pass.

We'll examine the way that PF evaluates rules and how ordering matters in a bit more detail in Chapter 3, in the context of a slightly longer rule set.

For any rule that has a `keep state` part, PF keeps information about the connection (including various counters and sequence numbers) as an entry in the *state table*. The state table is where PF keeps information about existing connections that have already matched a rule, and new packets that arrive are compared to existing state table entries to find a match first. Only when a packet does not match any existing state will PF move on to a full *rule set evaluation*, checking if the packet matches a rule in the loaded rule set. We can also instruct PF to act on state information in various ways, but in a simple case like this, our main goal is to allow return traffic for connections we initiate to return to us.

Note that on OpenBSD 4.1 and later, the default for `pass` rules is to keep state information,<sup>3</sup> and we no longer need to specify `keep state` explicitly in a simple case like this. This means the rule set could be written like this:

---

```
# minimal rule set, OpenBSD 4.1 onwards keeps state by default
block in all
pass out all
```

---

In fact, you could even leave out the `all` keyword here if you like.

The other BSDs have mostly caught up with this change by now, and for the rest of this book, we will stick to the newer style rules, with an occasional reminder in case you are using an older system.

It goes pretty much without saying that passing all traffic generated by a specific host implies that the host in question is, in fact, trustworthy. This is something you do only if this is a machine you know you can trust.

When you're ready to use this rule set, load it with the following:

---

```
$ sudo pfctl -ef /etc/pf.conf
```

---

The rule set should load without any error messages or warnings. On all but the slowest systems, you should be returned to the `$` prompt immediately.

<sup>3</sup> In fact, the new default corresponds to `keep state flags S/SA`, ensuring that only initial SYN packets during connection setup create state, eliminating some puzzling error scenarios. To filter statelessly, you can specify `no state` for the rules where you do not want to record or keep state information. On FreeBSD, OpenBSD 4.1-equivalent PF code was merged into version 7.0.

## **Testing the Rule Set**

It's always a good idea to test your rule sets to make sure that they work as expected. Proper testing will become essential once you move on to more complicated configurations.

To test the simple rule set, see if it can perform domain name resolution. For example, you could see if `$ host nostarch.com` returns information such as the IP address of the host *nostarch.com* and the hostnames of that domain's mail exchangers. Or just see if you can surf the Web. If you can connect to external websites by name, the rule set is able to perform domain name resolution. Basically, any service you try to access from your own system should work, and any service you try to access on your system from another machine should produce a `Connection refused` message.

## **Slightly Stricter: Using Lists and Macros for Readability**

The rule set in the previous section is an extremely simple one—probably too simplistic for practical use. But it's a useful starting point to build from to create a slightly more structured and complete setup. We'll start by denying all services and protocols, and then allow only those that we know that we need<sup>4</sup> using lists and macros for better readability and control.

A *list* is simply two or more objects of the same type that you can refer to in a rule set, such as this:

---

```
pass proto tcp to port { 22 80 443 }
```

---

Here, `{ 22 80 443 }` is a list.

A *macro* is a pure readability tool. If you have objects that you will refer to more than once in your configuration, such as an IP address for an important host, it could be useful to define a macro instead. For example, you might define this macro early in your rule set:

---

```
external_mail = 192.0.2.12
```

---

Then you could refer to that host as `$external_mail` later in the rule set:

---

```
pass proto tcp to $external_mail port 25
```

---

These two techniques have great potential for keeping your rule sets readable, and as such, they are important factors that contribute to the overall goal of keeping you in control of your network.

---

4. Why write the rule set to default deny? The short answer is that it gives you better control. The point of packet filtering is to take control, not to run catch-up with what the bad guys do. Marcus Ranum has written a very entertaining and informative article about this called “The Six Dumbest Ideas in Computer Security” ([http://www.ranum.com/security/computer\\_security/editorials/dumb/index.html](http://www.ranum.com/security/computer_security/editorials/dumb/index.html)).

## A Stricter Baseline Rule Set

Up to this point, we've been rather permissive with regard to any traffic we generate ourselves. A permissive rule set can be very useful while we check that basic connectivity is in place, or to check whether filtering is part of a problem we are seeing. Once the “Do we have connectivity?” phase is over, it's time to start tightening up to create a baseline that keeps us in control.

To begin, add the following rule to `/etc/pf.conf`:

---

```
block all
```

---

This rule is completely restrictive and will block all traffic in all directions. This is the initial baseline filtering rule that we'll use in all complete rule sets over the next few chapters. We basically start from zero, with a configuration where *nothing* is allowed to pass. Later on, we will add rules that cut our traffic some more slack, but we will do so incrementally and in a way that keeps us firmly in control.

Next, we'll define a few macros for later use in the rule set:

---

```
tcp_services = "{ ssh, smtp, domain, www, pop3, auth, https, pop3s }"  
udp_services = "{ domain }"
```

---

Here, you can see how the combination of lists and macros can be turned to our advantage. Macros can be lists, and as demonstrated in the example, PF understands rules that use the names of services as well as port numbers, as listed in your `/etc/services` file. We will take care to use all these elements and some further readability tricks as we tackle complex situations that require more elaborate rule sets.

Having defined these macros, we can use them in our rules, which we will now edit slightly to look like this:

---

```
block all  
pass out proto tcp to port $tcp_services  
pass proto udp to port $udp_services
```

---

**NOTE** *Be sure to add `keep state` to these rules if your system has a PF version older than OpenBSD 4.1.*

The strings `$tcp_services` and `$udp_services` are macro references. Macros that appear in a rule set are expanded in place when the rule set loads, and the running rule set will have the full lists inserted where the macros are referenced. Depending on the exact nature of the macros, they may cause single rules with macro references to expand into several rules. Even in a small rule set like this, the use of macros makes the rules easier to grasp and maintain. The amount of information that needs to appear in the rule set shrinks, and with sensible macro names, the logic becomes clearer. To follow the logic in a typical rule set, more often than not, we do not need to see full lists of IP addresses or port numbers in place of every macro reference.

From a practical rule set maintenance perspective, it is important to keep in mind which services to allow on which protocol in order to keep a comfortably tight regime. Keeping separate lists of allowed services according to protocol is likely to be useful in keeping your rule set both functional and readable.

### TCP VS. UDP

We have taken care to separate out UDP services from TCP services. Several services run primarily on well-known port numbers on either TCP or UDP, and a few alternate between using TCP and UDP according to specific conditions.

The two protocols are quite different in several respects. TCP is connection-oriented and reliable, a perfect candidate for stateful filtering. In contrast, UDP is stateless and connectionless, but PF creates and maintains data equivalent to state information for UDP traffic in order to ensure UDP return traffic is allowed back if it matches an existing state.

One common example where state information for UDP is useful is for handling name resolution. When you ask a name server to resolve a domain name to an IP address or to resolve an IP address back to a hostname, it's reasonable to assume that you want to receive the answer. Retaining state information or the functional equivalent about your UDP traffic makes this possible.

### ***Reloading the Rule Set and Looking for Errors***

After we've changed our *pf.conf* file, we need to load the new rules, as follows:

---

```
$ sudo pfctl -f /etc/pf.conf
```

---

If there are no syntax errors, pfctl should not display any messages during the rule load.

If you prefer to display verbose output, use the *-v* flag:

---

```
$ sudo pfctl -vf /etc/pf.conf
```

---

When you use verbose mode, pfctl should expand your macros into their separate rules before returning you to the command-line prompt, as follows:

---

```
$ sudo pfctl -vf /etc/pf.conf
tcp_services = "{ ssh, smtp, domain, www, pop3, auth, https, pop3s }"
udp_services = "{ domain }"
block drop all
pass out proto tcp from any to any port = ssh flags S/SA keep state
pass out proto tcp from any to any port = smtp flags S/SA keep state
pass out proto tcp from any to any port = domain flags S/SA keep state
pass out proto tcp from any to any port = www flags S/SA keep state
pass out proto tcp from any to any port = pop3 flags S/SA keep state
pass out proto tcp from any to any port = auth flags S/SA keep state
```

---

```
pass out proto tcp from any to any port = https flags S/SA keep state
pass out proto tcp from any to any port = pop3s flags S/SA keep state
pass proto udp from any to any port = domain keep state
$ _
```

---

Compare this output to the content of the `/etc/pf.conf` file you actually wrote. Our single TCP services rule is expanded into eight different ones: one for each service in the list. The single UDP rule takes care of only one service, and it expands from what we wrote to include the default options. Notice that the rules are displayed in full, with default values such as `flags S/SA keep state` applied in place of any options you do not specify explicitly. This is the configuration as it is actually loaded.

## ***Checking Your Rules***

If you have made extensive changes to your rule set, check them before attempting to load the rule set by using the following:

---

```
$ pfctl -nf /etc/pf.conf
```

---

The `-n` option tells PF to parse the rules only, without loading them—more or less as a dry run and to allow you to review and correct any errors. If `pfctl` finds any syntax errors in your rule set, it will exit with an error message that points to the line number where the error occurred.

Some firewall guides advise you to make sure that your old configuration is truly gone, or you will run into trouble—your firewall might be in some kind of intermediate state that does not match either the before or after state. That is simply not true when you’re using PF. The last valid rule set loaded is active until you either disable PF or load a new rule set. `pfctl` checks the syntax, and then loads your new rule set completely before switching over to the new one. Once a valid rule set has been loaded, there is no intermediate state with a partial rule set or no rules loaded. One consequence is that traffic that matches states that are valid in both the old and new rule set will not be disrupted.

Unless you have actually followed the advice from some of those old guides and *flushed* your existing rules (that *is* possible, using `pfctl -F all` or similar) before attempting to load a new one from your configuration file, the last valid configuration will remain loaded. In fact, flushing the rule set is rarely a good idea, since it effectively puts your packet filter in a `pass all` mode, and with a rather high risk of disrupting useful traffic while you are getting ready to load your rules.

## ***Testing the Changed Rule Set***

Once you have a rule set that `pfctl` loads without any errors, it’s time to see if the rules you have written behave as expected. Testing name resolution with a command such as `$ host nostarch.com` (as we did earlier) should still work,

but choose a domain you have not accessed recently (such as one for a political party you would not consider voting for) to be sure that you’re not pulling DNS information from the cache.

You should be able to surf the Web and use several mail-related services, but due to the nature of this updated rule set, attempts to access TCP services other than the ones defined (SSH, SMTP, and so on) on any remote system should fail. And, as with our simple rule set, your system should refuse all connections that do not match existing state table entries; only return traffic for connections initiated by this machine will be allowed in.

## Displaying Information About Your System

The tests you’ve performed on your rule sets should have shown that PF was running and that your rules are behaving as expected. There are several ways to keep track of what happens in your running system. One of the more straightforward ways of extracting information about PF is to use the already familiar `pfctl` program.

Once PF is enabled and running, the system updates various counters and statistics in response to network activity. To confirm that PF is actually running and to view statistics about its activity, you can use `pfctl -s`, followed by the type of information you want to display. A pretty long list of information types is available (see `man 8 pfctl` and look for the `-s` options). We will be getting back to some of those display options in Chapter 8, and go into further detail about some of the statistics they provide in Chapter 9 when we use the data to optimize the configuration we are building.

The following shows an example of just the top part of the output of `pfctl -s info` (taken from my home gateway). The high-level information that indicates the system actually passes traffic can be found in this upper part.

---

```
$ sudo pfctl -s info
Status: Enabled for 24 days 12:11:27          Debug: err

Interface Stats for nfe0           IPv4            IPv6
Bytes In                  43846385394          0
Bytes Out                 20023639992         64
Packets In
  Passed                  49380289          0
  Blocked                  49530           0
Packets Out
  Passed                  45701100           1
  Blocked                  1034            0

State Table                   Total          Rate
current entries                319
searches                      178598618        84.3/s
inserts                        4965347        2.3/s
removals                       4965028        2.3/s
```

---

The first line of the `pfctl` output indicates that PF is enabled and has been running for a little more than three weeks, which is equal to the time since I last performed a system upgrade that required a reboot.

The Interface Stats part of the display shows the bytes in and out handled by the interface. The next few items are likely to be more interesting in our context, showing the number of packets blocked or passed in each direction. This is where we find an early indication of whether the filtering rules we wrote are catching any traffic. In this case, either the rule set matches expected traffic well or we have fairly well-behaved users and guests, with the number of packets passed overwhelmingly larger than the number of packets blocked in both directions.

The next important indicator of a working system that is processing traffic is the block of State Table statistics. The state table current entries line shows that there are 319 active states or connections, while the state table has been searched (searches) for matches to existing states on average a little more than 84 times per second, for a total of just over 178 million times since the counters were reset. The inserts and removals counters show the number of times states have been created and removed, respectively. As expected, the number of insertions and removals differs by the number of currently active states, and the rate counters show that for the time since the counters were last reset, the rate of states created and removed matches exactly, up to the resolution of this display.

The information here is roughly in line with the statistics you should expect to see on a gateway for a small network configured for IPv4 only. There is no reason to be alarmed by the packet passed in the IPv6 column. OpenBSD comes with IPv6 built in. During network interface configuration, by default, the TCP/IP stack sends IPv6 neighbor solicitation requests for the link local address. In a normal IPv4-only configuration, only the first few packets actually pass, and by the time the PF rule set from `/etc/pf.conf` is fully loaded, IPv6 packets are blocked by the `block all` default rule. (In this example, they do not show up in `nfe0`'s statistics because IPv6 is tunneled over a different interface.)

## Looking Ahead

You should now have a machine that can communicate with other Internet-connected machines, using a very basic rule set that serves as a starting point for controlling your network traffic. As you progress through this book, you'll learn how to add rules that do various useful things. In Chapter 3, we will extend the configuration to act as a gateway for a small network. Serving the needs of several computers has some consequences, and we will look at how to let at least some ICMP and UDP traffic through, for your own troubleshooting needs if nothing else.

In Chapter 3, we'll also consider network services that have consequences for your security, like FTP. Using packet filtering intelligently to handle services that are demanding, security-wise, is a recurring theme in this book.



# 3

## INTO THE REAL WORLD



The previous chapter demonstrated the configuration for basic packet filtering on a single machine. In this chapter, we will build on that basic setup, but move into more conventional territory: the packet-filtering *gateway*. Although most of the items in this chapter are potentially useful in a single-machine setup, our main focus is to set up a gateway that forwards a selection of network traffic and handles common network services for a basic local network.

## A Simple Gateway

We will start with building what you probably associate with the term *firewall*: a machine that acts as a gateway for at least one other machine. In addition to forwarding packets between its various networks, this machine's mission will be to improve the signal-to-noise ratio in the network traffic it handles. That's where our PF configuration comes in.

But before diving into the practical configuration details, we need to dip into some theory and flesh out some concepts. Bear with me; this will end up saving you some headaches I've seen on mailing lists and newsgroups all too often.

### ***Keep It Simple: Avoid the Pitfalls of in, out, and on***

In the single-machine setup, life is relatively simple. Traffic you create should either pass out to the rest of the world or be blocked by your filtering rules, and you get to decide what you want to let in from elsewhere.

When you set up a gateway, your perspective changes. You go from the "It's me versus the network out there" mindset to "I'm the one who decides what to pass to or from all the networks I'm connected to." The machine has several, or at least two, network interfaces, each connected to a separate network, and its primary function (or at least the one we're interested in here) is to forward network traffic between networks.

It's very reasonable to think that if you want traffic to pass from the network connected to `re1` to hosts on the network connected to `re0`, you will need a rule like the following.<sup>1</sup>

---

```
pass in proto tcp on re1 from re1:network to re0:network \
    port $ports keep state
```

---

However, one of the most common and most complained-about mistakes in firewall configuration is not realizing that the `to` keyword does not in itself guarantee passage to the end point. The `to` keyword here means only that a packet or connection must have a destination address that matches those criteria in order to match the rule. The rule we just wrote lets the traffic pass `in` to just the gateway itself, `on` the specific interface named in the rule. To allow the packets in a bit further and to move on to the next network, we need a matching rule that says something like this:

---

```
pass out proto tcp on re0 from re1:network to re0:network \
    port $ports keep state
```

---

But please stop and take a moment to read those rules one more time. This last rule allows only packets with a destination in the network directly connected to `re0` to pass, and nothing else. If that's exactly what you want,

---

1. In fact, the `keep state` part denotes the default behavior and is redundant if you are working with a PF version taken from OpenBSD 4.1 or later. However, there is generally no need to remove the specification from existing rules you come across when upgrading from earlier versions. To ease transition, the examples in this book will make this distinction when needed.

fine. In other contexts, such rules are, while perfectly valid, more specific than the situation calls for. It's very easy to let yourself dive deeply into specific details and lose the higher-level view of the configuration's purpose, and maybe earn yourself a few extra rounds of debugging in the process.

If there are good reasons for writing very specific rules like the preceding ones, you probably already know that you need them and why. By the time you have finished this book (if not a bit earlier), you should be able to articulate the circumstances when more specific rules are needed. However, for the basic gateway configurations in this chapter, it is likely that you will want to write rules that are not interface-specific. In fact, in some cases, it is not useful to specify the direction either, and you would simply use a rule like the following to let your local network access the Internet.

---

```
pass proto tcp from re1:network to port $ports keep state
```

---

For simple setups, interface-bound `in` and `out` rules are likely to add more clutter to your rule sets than they are worth. For a busy network admin, a readable rule set is a safer one.

For the remainder of this book, with some exceptions, we will keep the rules as simple as possible for readability.

## ***Network Address Translation vs. IPv6***

Once we start handling traffic between separate networks, it's useful to look at how network addresses work and why you are likely to come across several different addressing schemes. The subject of network addresses has been a rich source of both confusion and buzzwords over the years. The underlying facts are sometimes hard to establish, unless you go to the source and wade through a series of RFCs. Over the next few paragraphs, I will make an effort to clear up some of the confusion.

For example, a widely held belief is that if you have an internal network that uses a totally different address range than the one assigned to the interface attached to the Internet, you're safe, and no one from the outside can get at your network resources. This belief is closely related to the idea that the IP address of your firewall in the local network must be either 192.168.0.1 or 10.0.0.1.

There is an element of truth in both notions, and those addresses are common defaults. But the real story is that it is possible to sniff one's way past network address translation (although PF offers some tricks that make that task a little harder).

The real reason we use a specific set of internal address ranges and a different set of addresses for unique external address ranges is not primarily due to security concerns, but because it was the easiest way to work around a design problem in the Internet protocols: a limited range of possible addresses.

In the 1980s, when the Internet protocols were formulated, most computers on the Internet (or ARPANET, as it was known at the time) were large machines with anything from several dozen to several thousand users each. At the time, a 32-bit address space with more than four billion addresses

seemed quite sufficient, but several factors have conspired to prove that assumption wrong. One factor is that the address-allocation process lead to a situation where the largest chunks of the available address space were already allocated before some of the world's more populous nations even connected to the Internet. The other, and perhaps more significant, factor was that by the early 1990s, the Internet was no longer a research project, but rather a commercially available resource with consumers and companies of all sizes consuming IP address space at an alarming rate.

The long-term solution was to redefine the Internet to use a larger address space. In 1998, the specification for IPv6, with 128 bits of address space for a total of  $2^{128}$  addresses, was published as RFC 2460. But while we were waiting for IPv6 to become generally available, we needed a stopgap solution. That solution came as a series of RFCs that specified how a gateway could forward traffic with IP addresses translated, so that a large local network would look like just one computer to the rest of the Internet. Certain previously unallocated IP address ranges were set aside for these private networks. These were free for anyone to use, on the condition that traffic in those ranges would not be allowed out on the Internet untranslated. Thus Network Address Translation (NAT) was born in the mid-1990s, and quickly became the default way to handle addressing in local networks.<sup>2</sup>

PF supports IPv6 as well as the various IPv4 address translation tricks. (In fact, the BSDs were among the earliest IPv6 adopters, thanks to the efforts of the KAME project.<sup>3</sup>) All systems that have PF also support both the IPv4 and the IPv6 address families. If your IPv4 network needs a NAT configuration, you can integrate the translation as needed in your PF rule set. In other words, if you are using a system that supports PF, you can rest assured that your IPv6 needs have been taken care of, at least on the operating system level.

The examples in this book use mainly IPv4 addresses and NAT where appropriate, but most of the material is equally relevant in networks that have implemented IPv6.

## **Final Preparations: Defining Your Local Network**

In Chapter 2, we set up a configuration for a single, stand-alone machine. We are about to extend that configuration to a gateway version, and it's useful to define a few more macros to help readability and to conceptually separate the local networks where you have a certain measure of control versus everything else. So how do you define your “local” network in PF terms?

Earlier in this chapter, you saw the *interface:network* notation. This is a nice piece of shorthand, but you can make your rule set even more readable and easier to maintain by taking the macro a bit further. For example, you could define a \$localnet macro as the network directly attached to your internal

---

2. RFC 1631, “The IP Network Address Translator (NAT),” dated May 1994, and RFC 1918, “Address Allocation for Private Internets,” dated February 1996, provide the details about NAT.

3. To quote the project home page at <http://www.kame.net/>, “The KAME project was a joint effort of six companies in Japan to provide a free stack of IPv6, IPsec, and Mobile IPv6 for BSD variants.” The main research and development activities were considered complete in March 2006, with only maintenance activity continuing, now that the important parts have been incorporated into the relevant systems.

interface (`re1:network` in our examples). Or you could change the definition of `$localnet` to an IP address/netmask notation to denote a network, such as `192.168.100.0/24` for a subnet of private IPv4 addresses, or `feco:dead:beef::/64` for an IPv6 range.

If your network environment requires it, you could define your `$localnet` as a list of networks. For example, a sensible `$localnet` definition combined with pass rules that use the macro, such as the following, could end up saving you a few headaches.

---

```
pass proto { tcp, udp } from $localnet to port $ports
```

---

We will stick to the convention of using macros such as `$localnet` for readability from here on.

## ***Setting Up a Gateway***

We will take the single-machine configuration we built from the ground up in the previous chapter as our starting point for building our packet-filtering gateway. We assume that the machine has acquired another network card (or you have set up a network connection from your local network to one or more other networks, via Ethernet, PPP, or other means).

In our context, it is not too interesting to look at the details of how the interfaces are configured. We just need to know that the interfaces are up and running.

For the following discussion and examples, only the interface names will differ between a PPP setup and an Ethernet one, and we will do our best to get rid of the actual interface names as quickly as possible.

First, since packet forwarding is off by default in all BSDs, we need to turn it on in order to let the machine forward the network traffic it receives on one interface to other networks via one or more separate interfaces. Initially, we will do this on the command line with a `sysctl` command, for traditional IPv4:

---

```
# sysctl net.inet.ip.forwarding=1
```

---

If we need to forward IPv6 traffic, we use this `sysctl` command:

---

```
# sysctl net.inet6.ip6.forwarding=1
```

---

This is fine for now. However, in order for this to work once you reboot the computer at some time in the future, you need to enter these settings into the relevant configuration files.

In OpenBSD and NetBSD, you do this by editing `/etc/sysctl.conf` and adding or changing the IP forwarding lines to look like this:

---

```
net.inet.ip.forwarding=1
net.inet6.ip6.forwarding=1
```

---

In FreeBSD, make the change by putting these lines in your */etc/rc.conf*:

---

```
gateway_enable="YES" #for ipv4
ipv6_gateway_enable="YES" #for ipv6
```

---

The net effect is identical; the FreeBSD *rc* script sets the two values via *sysctl* commands. However, a larger part of the FreeBSD configuration is centralized into the *rc.conf* file.

Now it's time to check whether all of the interfaces you intend to use are up and running. Use **ifconfig -a** or **ifconfig interface\_name** to find out.

The output of **ifconfig -a** on one of my systems looks like this:

---

```
$ ifconfig -a
lo0: flags=8049<UP,LOOPBACK,RUNNING,MULTICAST> mtu 33224
    groups: lo
        inet 127.0.0.1 netmask 0xffff0000
        inet6 ::1 prefixlen 128
            inet6 fe80::1%lo0 prefixlen 64 scopeid 0x5
xlo: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
    lladdr 00:60:97:83:4a:01
    groups: egress
    media: Ethernet autoselect (100baseTX full-duplex)
    status: active
        inet 194.54.107.18 netmask 0xffffffff broadcast 194.54.107.23
        inet6 fe80::260:97ff:fe83:4a01%xlo prefixlen 64 scopeid 0x1
fxp0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
    lladdr 00:30:05:03:fc:41
    media: Ethernet autoselect (100baseTX full-duplex)
    status: active
        inet 194.54.103.65 netmask 0xffffffffc0 broadcast 194.54.103.127
        inet6 fe80::230:5ff:fe03:fc41%fxp0 prefixlen 64 scopeid 0x2
pflog0: flags=141<UP,RUNNING,PROMISC> mtu 33224
enc0: flags=0<> mtu 1536
```

---

Your setup is most likely somewhat different. Here, the physical interfaces on the gateway are **xlo** and **fxp0**. The logical interfaces **lo0** (the loopback interface), **enc0** (the encapsulation interface for IPSEC use), and **pflog0** (the PF logging device) are probably on your system, too.

If you are on a dial-up connection, you probably use some variant of PPP for the Internet connection, and your external interface is the **tun0** pseudo-device. If your connection is via some sort of broadband connection, you may have an Ethernet interface to play with. However, if you are in the significant subset of ADSL users who use PPP over Ethernet (PPPoE), the correct external interface will be one of the pseudo-devices **tun0** or **pppoe0** (depending on whether you use userland **pppoe(8)** or kernel mode **pppoe(4)**), not the physical Ethernet interface.

Depending on your specific setup, you may need to do some other device-specific configuration for your interfaces. After you have that set up, you can move on to the TCP/IP level and deal with the packet-filtering configuration.

If you still intend to allow any traffic initiated by machines on the inside, your `/etc/pf.conf` for your initial gateway setup could look roughly like this:

---

```
ext_if = "re0" # macro for external interface - use tun0 or pppoe0 for PPPoE
int_if = "re1" # macro for internal interface
localnet = $int_if:network
# ext_if IP address could be dynamic, hence ($ext_if)
match out on $ext_if from $localnet nat-to ($ext_if)
block all
pass from { lo0, $localnet }
```

---

Note the use of macros to assign logical names to the network interfaces. Here, Realtek Ethernet cards are used, but this is the last time we will find this of any interest whatsoever in our context.

In truly simple setups like this one, we may not gain very much by using macros like these, but once the rule sets grow a little larger, you will learn to appreciate the readability they add.

Also note the `match` rule with `nat-to`. This is where you handle NAT from the nonroutable address inside your local network to the sole official address assigned to you. If your network uses official, routable addresses, you simply leave this line out of your configuration. The `match` rules, which were introduced in OpenBSD 4.6, can be used to apply actions when a connection matches the criteria without deciding whether a connection should be blocked or passed.

The parentheses surrounding the last part of the `match` rule (`($ext_if)`) are there to compensate for the possibility that the IP address of the external interface may be dynamically assigned. This detail will ensure that your network traffic runs without serious interruptions, even if the interface's IP address changes.

If your system runs a pre-OpenBSD 4.7 PF version, your first gateway rule set would look something like this:

---

```
ext_if = "re0" # macro for external interface - use tun0 or pppoe0 for PPPoE
int_if = "re1" # macro for internal interface
localnet = $int_if:network
# ext_if IP address could be dynamic, hence ($ext_if)
nat on $ext_if from $localnet to any -> ($ext_if)
block all
pass from { lo0, $localnet } to any keep state
```

---

The `nat` rule here handles the translation much like the `match` rule with `nat-to` in the previous example.

On the other hand, this rule set probably allows more traffic than what you actually want to pass out of your network. In one of the networks where I've done a bit of work, the main part of the rule set is based on a macro called `client_out`:

---

```
client_out = "{ ftp-data, ftp, ssh, domain, pop3, auth, nntp, http,
https, 446, cvspserver, 2628, 5999, 8000, 8080 }"
```

---

It has this pass rule:

---

```
pass inet proto tcp from $localnet to port $client_out
```

---

This may be a somewhat peculiar selection of ports, but it's exactly what my colleagues there needed at the time. Some of the numbered ports were needed for systems that were set up for specific purposes at other sites. Your needs probably differ at least in some details, but this should cover some of the more useful services.

Here's another pass rule that is useful to those who want the ability to administer machines from elsewhere:

---

```
pass in inet proto tcp to port ssh
```

---

Or use this form, if you prefer:

---

```
pass in inet proto tcp to $ext_if port ssh
```

---

When you leave out the `from` part entirely, the default used is `from any`, which is really quite permissive. It lets you log in from anywhere, which is great if you travel a lot and need SSH access from unknown locations around the world. If you're not all that mobile—say you haven't quite developed the taste for conferences in far-flung locations or you really want to leave your colleagues to fend for themselves while you're on vacation—you may want to tighten up with a `from` part that includes only the places where you and other administrators are likely to log in from for legitimate reasons.

Our very basic rule set is still not complete. Next, we need to make the name service work for our clients. We start with another macro at the start of our rule set.

---

```
udp_services = "{ domain, ntp }"
```

---

This is supplemented with a rule that passes the traffic we want through our firewall:

---

```
pass quick inet proto { tcp, udp } to port $udp_services
```

---

Note the `quick` keyword in this rule. We have started writing rule sets that consist of several rules, and it is time to revisit the relationships and interactions between them.

As noted in the previous chapter, the rules are evaluated from top to bottom, in the sequence they are written in the configuration file. For each packet or connection evaluated by PF, *the last matching rule* in the rule set is the one that is applied.

The `quick` keyword offers an escape from the ordinary sequence. When a packet matches a `quick` rule, the packet is treated according to the present rule. The rule processing stops without considering any further rules that

might have matched the packet. As your rule sets grow longer and more complicated, you will find this quite handy. For example, it's useful when you need a few isolated exceptions to your general rules.

This quick rule also takes care of NTP, which is used for time synchronization. Common to both the name service and time synchronization protocols is that they may, under certain circumstances, communicate alternately over TCP and UDP.

## **Testing Your Rule Set**

You may not have gotten around to writing that formal test suite for your rule sets just yet, but there is every reason to test that your configuration works as expected.

The same basic tests in the stand-alone example from the previous chapter still apply. But now you need to test from the other hosts in your network as well as from your packet-filtering gateway. For each of the services you specified in your `pass` rules, test that machines in your local network get meaningful results. From any machine in your local network, enter a command like this:

---

```
$ host nostarch.com
```

---

It should return exactly the same results as when you tested the stand-alone rule set in the previous chapter, and traffic for the services you have specified should pass.<sup>4</sup>

You may not think it's necessary, but it does not hurt to check that the rule set works as expected from outside your gateway as well. If you have done exactly what this chapter says so far, it should not be possible to contact machines in your local network from the outside.

### **WHY ONLY IP ADDRESSES— NO HOSTNAMES OR DOMAIN NAMES?**

Looking at the examples up to this point, you have probably noticed that the rule sets all have macros that expand into IP addresses or address ranges, but never hostnames or domain names. You're probably wondering why. After all, you've seen that PF lets you use service names in your rule set, so why not include hostnames or domain names?

The answer is that if you used domain names and hostnames in your rule set, the rule set would be valid only after the name service is running and accessible. In the default configuration, PF is loaded before any network services are running. This means that if you want to use domain names and hostnames in your PF configuration, you will need to change the system's startup sequence (by editing `/etc/rc.local`, perhaps) to load the name service-dependent rule set only after the name service is available.

4. Unless, of course, the information changed in the meantime. Some sysadmins are fond of practical jokes, but most of the time changes in DNS zone information are due to real-world needs in that particular organization or network.

## That Sad Old FTP Thing

The short list of real-life TCP ports we looked at a few moments back contained, among other things, FTP, the classic file transfer protocol. FTP is a relic of the early Internet, when experiments were the norm, and security was not really on the horizon in any modern sense. FTP actually predates TCP/IP,<sup>5</sup> and it is possible to track the protocol's development through more than 50 RFCs. After more than 30 years, FTP is both a sad old thing and a problem child—emphatically so for anyone trying to combine FTP and firewalls. FTP is an old and weird protocol, with a lot to dislike. Here are the most common points against it:

- Passwords are transferred in the clear.
- The protocol demands the use of at least two TCP connections (control and data) on separate ports.
- When a session is established, data is communicated via ports selected at random.

All of these points make for challenges security-wise, even before considering any potential weaknesses in client or server software that may lead to security issues. As any network graybeard will tell you, these things tend to crop up when you need them the least.

Under any circumstances, other, more modern and more secure options for file transfer exist, such as SFTP and SCP, which feature both authentication and data transfer via encrypted connections. Competent IT professionals should have a preference for some other form of file transfer than FTP.

Regardless of our professionalism and preferences, we sometimes must handle things we would prefer not to use at all. In the case of FTP through firewalls, we can combat problems by redirecting the traffic to a small program that is written specifically for this purpose. The upside for us is that handling FTP offers us a chance to look at *redirection*.

The easiest way to handle FTP is to have PF redirect the traffic for that service to an external application that acts as a *proxy* for the service. The proxy then interacts with your packet filter through a well-defined interface.

### If We Must: *ftp-proxy* with Redirection

Enabling FTP transfers through your gateway is amazingly simple, thanks to the FTP proxy program included in the OpenBSD base system. The program is called—you guessed it—*ftp-proxy*.

FTP being what it is, the proxy needs to dynamically insert rules in your rule set. *ftp-proxy* interacts with your configuration via an *anchor* (a named rule set section) where the proxy inserts and deletes the rules it constructs to handle your FTP traffic.

---

5. The earliest RFC describing the FTP protocol is RFC 114, dated April 10, 1971. The switch to TCP/IP happened with FTP version 5 as defined in RFCs 765 and 775, dated June and December 1980, respectively.

To enable `ftp-proxy`, you need to add this line to your `/etc/rc.conf.local` file:

---

```
ftpproxy_flags=""
```

---

You can start the proxy manually by running `/usr/sbin/ftp-proxy` if you like, and you may want to do this in order to check that the changes to the PF configuration you are about to make have the intended effect.

For a basic configuration, you need to add only three elements to your `/etc/pf.conf`: the anchor and two pass rules. The anchor declaration looks like this:

---

```
anchor "ftp-proxy/*"
```

---

In pre-OpenBSD 4.7 versions, two anchor declarations were needed:

---

```
nat-anchor "ftp-proxy/*"  
rdr-anchor "ftp-proxy/*"
```

---

The proxy will insert the rules it generates for the FTP sessions here. Then you also need a `pass` rule to let FTP traffic in to the proxy:

---

```
pass in quick proto tcp to port ftp rdr-to 127.0.0.1 port 8021
```

---

Note the `rdr-to` part. This redirects the traffic to the local port where the proxy listens.

Here is the pre-OpenBSD 4.7 version of this rule:

---

```
rdr pass on $int_if proto tcp from any to any port ftp -> 127.0.0.1 port 8021
```

---

Finally, add a `pass` rule to let the packets pass from the proxy to the rest of the world:

---

```
pass out proto tcp from $proxy to any port ftp
```

---

`$proxy` expands to the address to which the proxy daemon is bound. Reload your PF configuration:

---

```
$ sudo pfctl -f /etc/pf.conf
```

---

Before you know it, your users will thank you for making FTP work.

This example covers a basic setup where the clients in your local network need to contact FTP servers elsewhere. The basic configuration here should work well with most combinations of FTP clients and servers.

You can change the proxy's behavior in various ways by adding options to the `ftpproxy_flags=` line. You may bump into clients or servers with specific quirks that you need to compensate for in your configuration, or you may want to integrate the proxy in your setup in specific ways, such as assigning FTP traffic to a specific queue. For these and other finer points of `ftp-proxy` configuration, your best bet is to start by studying the man page.

If you are interested in ways to run an FTP server protected by PF and `ftp-proxy`, you could look into running a separate `ftp-proxy` in reverse mode (using the `-R` option), on a separate port with its own redirecting `pass` rule.

**NOTE** *If your PF version predates the ones described here, please look up the first edition of this book and the documentation for your operating system for information on how to use some earlier FTP proxies.*

## Making Your Network Troubleshooting Friendly

Making your network troubleshooting friendly is a potentially large subject. Generally, the debugging or troubleshooting friendliness of your TCP/IP network depends on how you treat the Internet protocol that was designed specifically with debugging in mind: ICMP.

ICMP is the protocol for sending and receiving *control messages* between hosts and gateways, mainly to provide feedback to a sender about any unusual or difficult conditions en route to the target host.

There is a lot of ICMP traffic, which usually happens in the background while you are surfing the Web, reading mail, or transferring files. Routers (you are aware that you are building one, right?) use ICMP to negotiate packet sizes and other transmission parameters in a process often referred to as *path MTU discovery*.

You may have heard admins referring to ICMP as either “evil,” or, if their understanding runs a little deeper, “a necessary evil.” The reason for this attitude is purely historical. A few years back, it was discovered that the networking stacks of several operating systems contained code that could make the machine crash if it was sent a sufficiently large ICMP request.

One of the companies that was hit hard by this was Microsoft, and you can find a lot of material on the *ping of death* bug by using your favorite search engine. However, this all happened in the second half of the 1990s, and all modern operating systems have thoroughly sanitized their network code since then (at least, that’s what we are lead to believe).

One of the early work-arounds was to simply block either ICMP echo (ping) requests or even all ICMP traffic. Now these rule sets have been around for roughly 10 years, and the people who put them there are still scared. There is most likely little or no reason to worry about destructive ICMP traffic anymore, but here we will look at how to manage just what ICMP traffic passes to or from your network.

## ***Do We Let It All Through?***

The obvious question becomes, “If ICMP is such a good and useful thing, shouldn’t we let it all through, all the time?” The answer is that it depends.

Letting diagnostic traffic pass unconditionally makes debugging easier, of course, but it also makes it relatively easy for others to extract information about your network. So, a rule like the following might not be optimal if you want to cloak the internal workings of your network:

---

```
pass inet proto icmp
```

---

In all fairness, it should also be said that you might find some ICMP traffic quite harmlessly riding piggyback on your `keep_state` rules.

## ***The Easy Way Out: The Buck Stops Here***

The easiest solution could very well be to allow all ICMP traffic from your local network through, and let probes from elsewhere stop at your gateway:

---

```
pass inet proto icmp icmp-type $icmp_types from $localnet
pass inet proto icmp icmp-type $icmp_types to $ext_if
```

---

Stopping probes at the gateway might be an attractive option anyway, but let’s look at a few other options that will demonstrate some of PF’s flexibility.

## ***Letting ping Through***

The rule set we have developed so far has one clear disadvantage: Common troubleshooting commands such as `ping` and `traceroute` will not work. That may not matter too much to your users, and since it was the `ping` command that scared people into filtering or blocking ICMP traffic in the first place, there are apparently some people who feel we are better off without it. However, if you are in my perceived target audience, you will be rather fond of having those troubleshooting tools available. And with a couple of small additions to the rule set, they will be.

`ping` uses ICMP, and in order to keep our rule set tidy, we start by defining another macro:

---

```
icmp_types = "echoreq"
```

---

Then we add a rule that uses the definition:

---

```
pass inet proto icmp icmp-type $icmp_types
```

---

If you need more or other types of ICMP packets to go through, you can expand `icmp_types` to a list of those packet types you want to allow.

## **Helping traceroute**

traceroute is another command that is quite useful when your users claim that the Internet isn't working. By default, Unix traceroute uses UDP connections according to a set formula based on destination. The following rule works with the traceroute command on all forms of Unix I've had access to, including GNU/Linux:

---

```
# allow out the default range for traceroute(8):
# "base+nhops*nqueries-1" (33434+64*3-1)
pass out on $ext_if inet proto udp to port 33433 >< 33626
```

---

This gives you a first taste of what port ranges look like. They are quite useful in some contexts.

Experience so far indicates that traceroute implementations on other operating systems work roughly the same. One notable exception is Microsoft Windows. On that platform, the *tracert.exe* program uses ICMP echo requests for this purpose. So if you want to let Windows traceroutes through, you need only the first rule, much like letting ping through. The Unix traceroute program can be instructed to use other protocols as well, and will behave remarkably like its Microsoft counterpart if you use its *-I* command-line option. You can check the traceroute man page (or its source code, for that matter) for all the details.

This solution is based on a sample rule I found in an *openbsd-misc* post. I've found that list, and the searchable list archives (accessible among other places from <http://marc.info/>), to be a very valuable resource whenever you need OpenBSD or PF-related information.

## **Path MTU Discovery**

The last bit I will remind you about when it comes to troubleshooting is the *path MTU discovery*. Internet protocols are designed to be device-independent, and one consequence of device independence is that you cannot always predict reliably what the optimal packet size is for a given connection. The main constraint on your packet size is called the *maximum transmission unit*, or *MTU*, which sets the upper limit on the packet size for an interface. The *ifconfig* command will show you the MTU for your network interfaces.

Modern TCP/IP implementations expect to be able to determine the correct packet size for a connection through a process that simply involves sending packets of varying sizes within the MTU of the local link with the “do not fragment” flag set. If a packet then exceeds the MTU somewhere along the way to the destination, the host with the lower MTU will return an ICMP packet indicating “type 3, code 4,” when the local upper limit has been reached. Now, you don’t need to dive for the RFCs right away. Type 3 means *destination unreachable*, and code 4 is short for *fragmentation needed, but the do not fragment*

*flag* is set. So if your connections to other networks, which may have MTUs that differ from your own, seem suboptimal, you could try changing your list of ICMP types slightly to let the destination-unreachable packets through:

---

```
icmp_types = "{ echoreq, unreach }"
```

---

As you can see, this means you do not need to change the pass rule itself:

---

```
pass inet proto icmp icmp-type $icmp_types
```

---

Now I'll let you in on a little secret: In almost all cases, these rules are not necessary for purposes of path MTU discovery (but they don't hurt either). However, even though the default PF keep state behavior takes care of most of the ICMP traffic you will need, PF does let you filter on all variations of ICMP types and codes. If you want to delve into more detail, the list of possible types and codes are documented in the `icmp(4)` and `icmp6(4)` man pages. The background information is available in the RFCs.<sup>6</sup>

## Tables Make Your Life Easier

By now, you may be thinking that this setup gets awfully static and rigid. There will, after all, be some kinds of data relevant to filtering and redirection at a given time, but they do not deserve to be put into a configuration file! Quite right, and PF offers mechanisms for handling those situations.

Tables are one such feature. They are useful as lists of IP addresses that can be manipulated without reloading the entire rule set and also when fast lookups are desirable.

Table names are always enclosed in `< >`, like this:

---

```
table <clients> persist { 192.168.2.0/24, !192.168.2.5 }
```

---

Here, the network `192.168.2.0/24` is part of the table with one exception: The address `192.168.2.5` is excluded using the `!` operator (logical NOT). The keyword `persist` makes sure the table itself will exist, even if no rules currently refer to it.

It is also possible to load tables from files where each item is on a separate line, such as the file `/etc/clients`:

---

```
192.168.2.0/24
!192.168.2.5
```

---

This, in turn, is used to initialize the table in `/etc/pf.conf`:

---

```
table <clients> persist file "/etc/clients"
```

---

<sup>6</sup>. The main RFCs describing ICMP and some related techniques are 792, 950, 1191, 1256, 2521, and 2765. ICMP updates for IPv6 are in RFC 1885, RFC 2463, and RFC 2466. These documents are available in a number of places on the Net, such as <http://www.ietf.org/> and <http://www.faqs.org/>, and probably also via your package system.

So, for example, you can change one of our earlier rules to read like this to manage outgoing traffic from your client computers:

---

```
pass inet proto tcp from <clients> to any port $client_out
```

---

With this in hand, you can manipulate the table's contents live, like this:

---

```
$ sudo pfctl -t clients -T add 192.168.1/16
```

---

Note that this changes the in-memory copy of the table only, meaning that the change will not survive a power failure or reboot, unless you arrange to store your changes.

You might opt to maintain the on-disk copy of the table using a cron job that dumps the table content to disk at regular intervals, using a command such as the following:

---

```
$ sudo pfctl -t clients -T show >/etc/clients
```

---

Alternatively, you could edit the */etc/clients* file and replace the in-memory table contents with the file data:

---

```
$ sudo pfctl -t clients -T replace -f /etc/clients
```

---

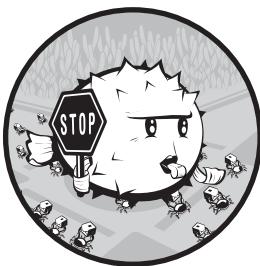
For operations you will be performing frequently, sooner or later, you will end up writing shell scripts. It is likely that routine operations on tables such as inserting or removing items or replacing table contents will be part of your housekeeping scripts in the near future.

One common example is to enforce network access restrictions via cron jobs that replace the contents of the tables referenced as `from` addresses in the `pass` rules at specific times. In some networks, you may even need different access rules for different days of the week. The only real limitations lie in your own needs and your creativity.

We will be returning to some handy uses of tables frequently over the next chapters, and we will look at a few programs that interact with tables in useful ways.

# 4

## WIRELESS NETWORKS MADE EASY



It is rather tempting to say that on BSD, and OpenBSD in particular, there's no need to "make wireless networking easy," because it already is. Getting a wireless network running is not very different from getting a wired one up and running, but there are some issues that turn up simply because we are dealing with radio waves and not wires. We will look briefly at some of the issues before moving on to the practical steps involved in creating a usable setup.

Once we have covered the basics of getting a wireless network up and running, we'll turn to some of the options for making your wireless network more interesting and harder to break.

## A Little IEEE 802.11 Background

Setting up any network interface, in principle, is a two-step process: you establish a link, and then you move on to configuring the interface for TCP/IP traffic.

In the case of wired Ethernet-type interfaces, establishing the link usually consists of plugging in a cable and seeing the link indicator light up. However, some interfaces require extra steps. Networking over dial-up connections, for example, requires telephony steps such as dialing a number to get a carrier signal.

In the case of IEEE 802.11-style wireless networks, getting the carrier signal involves quite a few steps at the lowest level. First, you need to select the proper channel in the assigned frequency spectrum. Once you find a signal, you need to set a few link-level network identification parameters. Finally, if the station you want to link to uses some form of link-level encryption, you need to set the correct kind and probably negotiate some additional parameters.

Fortunately, on BSD systems, all configuration of wireless network devices happens via `ifconfig` commands and options, as with any other network interface.<sup>1</sup> Still, since we are introducing wireless networks here, we need to look at security at various levels in the networking stack from this new perspective.

There are basically three kinds of popular and simple IEEE 802.11 security mechanisms, and we will discuss them briefly over the next sections.

**NOTE** *For a more complete overview of issues surrounding security in wireless networks, see Professor Kjell Jørgen Hole’s articles and slides at <http://www.kjhole.com/> and <http://www.kjhole.com/Standards/WiFi/WiFiDownloads.html>. For fresh developments in the Wi-Fi field, the Wi-Fi Net News site ([http://wifinetnews.com/archives/cat\\_security.html](http://wifinetnews.com/archives/cat_security.html)) and The Unofficial 802.11 Security Web Page (<http://www.drizzle.com/~aboba/IEEE/>) are highly recommended.*

### MAC Address Filtering

The short version of the story about PF and MAC address filtering is that we don’t do it.

A number of consumer-grade, off-the-shelf wireless access points offer MAC address filtering, but contrary to common belief, they don’t really add much security. The marketing succeeds largely because most consumers are unaware that it’s possible to change the MAC address of essentially any wireless network adapter on the market today.<sup>2</sup>

---

1. On some systems, the older, device-specific programs such as `wicontrol` and `anicontrol` are still around, but for the most part, they are deprecated and in the process of being replaced with `ifconfig` functionality. On OpenBSD, the consolidation into `ifconfig` has been completed.

2. A quick man page lookup on OpenBSD will tell you that the command to change the MAC address for the interface `xum0` is simply `ifconfig xum0 lladdr 00:ba:ad:fo:od:11`.

**NOTE** If you really want to try MAC address filtering, you could look into using the bridge(4) facility and the MAC filtering features offered by `brconfig(8)` (on OpenBSD 4.6 and earlier) or the bridge-related rule options in `ifconfig(8)` (OpenBSD 4.7 and later). We'll look at bridges and some of the more useful ways to use them with packet filtering in Chapter 5.

## WEP

One consequence of using radio waves instead of wires to move data is that it is comparatively easy for outsiders to capture data in transit over radio waves. The designers of the 802.11 family of wireless network standards seem to have been aware of this fact, and they came up with a solution that they went on to market under the name *Wired Equivalent Privacy*, or **WEP**.

Unfortunately, the WEP designers came up with their wired equivalent encryption without actually reading up on recent research or consulting active researchers in the field. So, the link-level encryption scheme they recommended is considered a pretty primitive homebrew among cryptography professionals. It was no great surprise when WEP encryption was reverse-engineered and cracked within a few months after the first products were released.

Even though you can download free tools to descramble WEP-encoded traffic in a matter of minutes, for a variety of reasons, WEP is still widely supported and used. Essentially, all IEEE 802.11 equipment available today has support for at least WEP, and a surprising number offer MAC address filtering, too.

You should consider network traffic protected only by WEP to be just marginally more secure than data broadcast in the clear. Then again, the token effort needed to crack into a WEP network may be sufficient to deter lazy and unsophisticated attackers.

## WPA

It dawned on the 802.11 designers fairly quickly that their WEP system was not quite what it was cracked up to be, so they came up with a revised and slightly more comprehensive solution called *Wi-Fi Protected Access*, or **WPA**.

WPA looks better than WEP, at least on paper, but the specification is complicated enough that its widespread implementation was delayed. In addition, WPA has attracted its share of criticism over design issues and bugs. Combined with the familiar issues of access to documentation and hardware, free software support varies. Most free systems have WPA support, but you may find that it is not available for all devices. If your project specification includes WPA, look carefully at your operating system and driver documentation.

And, of course, it goes almost without saying that you will need further security measures, such as SSH or SSL encryption, to maintain any significant level of confidentiality for your data stream.

## **The Right Hardware for the Task**

Picking the right hardware is not necessarily a daunting task. On a BSD system, the following simple command is all you need to enter to see a listing of all manual pages with the word *wireless* in their subject lines.<sup>3</sup>

---

```
$ apropos wireless
```

---

Even on a freshly installed system, this command will give you a complete list of all wireless network drivers available in the operating system.

The next step is to read the driver manual pages and compare the lists of compatible devices with what is available as parts or built into the systems you are considering. Take some time to think through your specific requirements. For test purposes, low-end `rum` or `ural` USB dongles will work. Later, when you are about to build a more permanent infrastructure, you may want to look into higher-end gear. You may also want to read Appendix B of this book.

## **Setting Up a Simple Wireless Network**

For our first wireless network, it makes sense to use the basic gateway configuration from the previous chapter as our starting point. In your network design, it is likely that the wireless network is not directly attached to the Internet at large, but the wireless network will require a gateway of some sort. For that reason, it makes sense to reuse the working gateway setup for this wireless access point, with some minor modifications introduced over the next few paragraphs. After all, doing so is more convenient than starting a new configuration from scratch.

**NOTE** *We are in infrastructure-building mode here, and will be setting up the access point first. If you prefer to look at the client setup first, see “The Client Side” on page 51.*

The first step is to make sure you have a supported card and check that the driver loads and initializes the card properly. The boot-time system messages scroll by on the console, but they are also recorded in the file `/var/run/dmesg.boot`. You can view the file itself or use the `dmesg` command to see these messages. With a successfully configured PCI card, you should see something like this:

---

```
ralo at pci1 dev 10 function 0 "Ralink RT2561S" rev 0x00: apic 2 int 11 (irq 11), address 00:25:9c:72:cf:60
ralo: MAC/BBP RT2561C, RF RT2527
```

---

If the interface you want to configure is a hot-pluggable type, such as a USB or PC Card device, you can see the kernel messages by viewing the `/var/log/messages` file; for example, by running `tail -f` on the file before you plug in the device.

---

<sup>3</sup>. In addition, it is possible to look up man pages on the Web. Check <http://www.openbsd.org/> and the other project websites. They offer keyword-based man page searching.

Next, you need to configure the interface, first to enable the link, and finally to configure the system for TCP/IP. You can do this from the command line, like this:

---

```
$ sudo ifconfig ralo up mediaopt hostap mode 11g chan 1 nwid unwiredbsd nwkey 0x1deadbeef9
```

---

This command does several things at once. It configures the `ralo` interface, enables the interface with the `up` parameter, and specifies that the interface is an access point for a wireless network with `mediaopt hostap`. Then it explicitly sets the operating mode to `11g` and the channel to `11`. Finally, it uses the `nwid` parameter to set the network name to `unwiredbsd`, with the WEP key (`nwkey`) set to the hexadecimal string `0x1deadbeef9`.

Use `ifconfig` to check that the command successfully configured the interface:

---

```
$ ifconfig ralo
ralo: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 00:25:9c:72:cf:60
      priority: 4
      groups: wlan
      media: IEEE802.11 autoselect mode 11g hostap
      status: active
      ieee80211: nwid unwiredbsd chan 1 bssid 00:25:9c:72:cf:60 nwkey <not displayed> 100dBm
      inet6 fe80::225:9cff:fe72:cf60%ralo prefixlen 64 scopeid 0x2
```

---

Note the contents of the `media` and `ieee80211` lines. The information displayed here should match what you entered on the `ifconfig` command line.

With the link part of your wireless network operational, you can assign an IP address to the interface:

---

```
$ sudo ifconfig ralo 10.50.90.1
```

---

On OpenBSD, you can combine both steps into one by creating a `/etc/hostname.ralo` file roughly like this:

---

```
up mediaopt hostap mode 11g chan 1 nwid unwiredbsd nwkey 0x1deadbeef9
inet 10.50.90.1
```

---

Then run `sh /etc/netstart ralo` (as root), or wait patiently for your next boot to complete.

Notice that the preceding configuration is divided over two lines. The first line generates an `ifconfig` command that sets up the interface with the correct parameters for the physical wireless network. The second line generates the command that sets the IP address after the first command completes. Because this is our access point, we set the channel explicitly, and we enable weak WEP encryption by setting the `nwkey` parameter.

On NetBSD, you can normally combine all of these parameters in one *rc.conf* setting:

---

```
ifconfig_ralo="mediaopt hostap mode 11g chan 1 nwid unwiredbsd nwkey 0x1deadbeef inet 10.50.90.1"
```

---

FreeBSD 8 and newer take a slightly different approach, tying wireless network devices to the unified *wlan(4)* driver. Depending on your kernel configuration, you may need to add the relevant module load lines to */boot/loader.conf*. On one of my test systems, */boot/loader.conf* looks like this:

---

```
if_rum_load="YES"
wlan_scan_ap_load="YES"
wlan_scan_sta_load="YES"
wlan_wep_load="YES"
wlan_ccmp_load="YES"
wlan_tkip_load="YES"
```

---

With the relevant modules loaded, setup is a multicommand affair, best handled by a *start\_if.iffile* for your wireless network. Here is an example of an */etc/start\_if.rum0* file for a WEP access point on FreeBSD 8:

---

```
wlans_rum0="wlano"
create_args_wlano="wlandev rum0 wlanmode hostap"
ifconfig_wlano="inet 10.50.90.1 netmask 255.255.255.0 ssid unwiredbsd \
wepmode on wepkey 0x1deadbeef9 mode 11g"
```

---

After a successful configuration, your *ifconfig* output should show both the physical interface and the *wlan* interface up and running:

---

```
rum0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 2290
      ether 00:24:1d:9a:bf:67
      media: IEEE 802.11 Wireless Ethernet autoselect mode 11g <hostap>
      status: running
wlan0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 1500
      ether 00:24:1d:9a:bf:67
      inet 10.50.90.1 netmask 0xffffffff broadcast 10.50.90.255
      media: IEEE 802.11 Wireless Ethernet autoselect mode 11g <hostap>
      status: running
      ssid unwiredbsd channel 6 (2437 Mhz 11g) bssid 00:24:1d:9a:bf:67
      country US authmode OPEN privacy ON deftxkey UNDEF wepkey 1:40-bit
      txpower 0 scanvalid 60 protmode CTS dtimperiod 1 -dfs
```

---

The line *status: running* means that you are up and running, at least on the link level.

**NOTE** Be sure to check the most up-to-date *ifconfig* man page for other options that may be more appropriate for your configuration.

## An OpenBSD WPA Access Point

WPA support was introduced in OpenBSD 4.4, with extensions to most wireless network drivers and a new utility called `wpa-psk(8)` to handle WPA keying operations.

**NOTE** *There may still be wireless network drivers that do not have WPA support, so check the driver's man page to see if WPA is supported before you try to configure your network to use it. You can combine 802.1x key management with an external authentication server for "enterprise" mode, but we will stick to the simpler preshared key setup for our purposes.*

The procedure for setting up an access point with WPA is quite similar to the one we followed for WEP. To generate a shared WPA key, you need to run the `wpa-psk` utility. If we reuse the WEP key from the earlier examples as the cleartext for our WPA passphrase, we could generate our key like this:

---

```
$ wpa-psk unwiredbsd 0x1deadbeef9  
0x31db31f2291f1ddf3ded3ca463a7dd5c0cd77a814f1d8e6c64990bfcb287b202
```

---

You could copy this value into the `ifconfig` command or *hostname.iffile*, or make `ifconfig` read the output of the `wpa-psk` call directly. Putting the cleartext into the configuration file will also make it slightly more readable. For a WPA setup with a preshared key (sometimes referred to as a *network password*), you would typically write a *hostname.iffile* like this:

---

```
up media autoselect mediaopt hostap mode 11g chan 1 nwid unwiredbsd \  
        wpa wpapsk `wpa-psk unwiredbsd 0x1deadbeef9`  
inet 10.50.90.1
```

---

If you are already running the WEP setup described earlier, disable those settings with the following:

---

```
$ sudo ifconfig ralo -nwid -nwkey
```

---

Then enable the new settings with this command:

---

```
$ sudo sh /etc/netstart ralo
```

---

You can then check that the access point is up and running with `ifconfig`:

---

```
$ ifconfig ralo  
ralo: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500  
      lladdr 00:25:9c:72:cf:60  
      priority: 4  
      groups: wlan
```

---

```
media: IEEE802.11 autoselect mode 11g hostap
status: active
ieee80211: nwid unwiredbsd chan 1 bssid 00:25:9c:72:cf:60 wpapsk <not displayed>
wpa_protos wpa1,wpa2 wpaakms psk wpaciphers tkip,ccmp wpagroupcipher tkip 100dBm
inet6 fe80::225:9cff:fe72:cf60%ralo prefixlen 64 scopeid 0x2
inet 10.50.90.1 netmask 0xffff0000 broadcast 10.255.255.255
```

---

Note the status: active indication and that the WPA options we did not set explicitly are shown with their sensible default values.

## A FreeBSD WPA Access Point

Moving from the WEP access point we configured earlier to a somewhat safer WPA setup is rather straightforward. WPA support on FreeBSD comes in the form of hostapd (a program that is somewhat similar to OpenBSD's hostapd but not identical). We start by editing the */etc/start\_if.rum0* file to remove the authentication information. The edited file should look something like this:

---

```
wlans_rum0="wlan0"
create_args_wlan0="wlandev rum0 wlanmode hostap"
ifconfig_wlan0="inet 10.50.90.1 netmask 255.255.255.0 ssid unwiredbsd mode 11g"
```

---

Next, we add the enable line for hostapd in */etc/rc.conf*.

---

```
hostapd_enable="YES"
```

---

And finally, hostapd will need some configuration of its own, in */etc/hostapd.conf*:

---

```
interface=wlan0
debug=1
ctrl_interface=/var/run/hostapd
ctrl_interface_group=wheel
ssid=unwiredbsd
wpa=1
wpa_passphrase=0x1deadbeef9
wpa_key_mgmt=WPA-PSK
wpa_pairwise=CCMP TKIP
```

---

Here, the interface specification is rather self-explanatory, while the debug value is set to produce minimal messages. The range is 0 through 4, where 0 is no debug messages at all. You should not need to change the ctrl\_interface\* settings unless you are developing hostapd. The first of the next five lines sets the network identifier. The subsequent lines enable WPA and set the pass-phrase. The final two lines specify accepted key-management algorithms and encryption schemes. (For the finer details and updates, see the *hostapd(8)* and *hostapd.conf(5)* man pages.)

After a successful configuration (running `sudo /etc/rc.d/hostapd forcestart` comes to mind), ifconfig should produce output about the two interfaces similar to this:

---

```
rum0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 2290
      ether 00:24:1d:9a:bf:67
      media: IEEE 802.11 Wireless Ethernet autoselect mode 11g <hostap>
      status: running
wlan0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 1500
      ether 00:24:1d:9a:bf:67
      inet 10.50.90.1 netmask 0xffffffff broadcast 10.50.90.255
      media: IEEE 802.11 Wireless Ethernet autoselect mode 11g <hostap>
      status: running
      ssid unwiredbsd channel 6 (2437 Mhz 11g) bssid 00:24:1d:9a:bf:67
      country US authmode WPA privacy MIXED deftxkey 2 TKIP 2:128-bit
      txpower 0 scanvalid 60 protmode CTS dtimperiod 1 -dfs
```

---

The line `status: running` means that you are up and running, at least on the link level.

### ***The Access Point's PF Rule Set***

With the interfaces configured, it's time to start configuring the access point as a packet-filtering gateway. You can start by copying the basic gateway setup from Chapter 3. Enable gatewaying via the appropriate entries in the access point's `sysctl.conf` or `rc.conf` file, and then copy across the `pf.conf` file. Depending on the parts of the previous chapter that were most useful to you, the `pf.conf` file may look something like this:

---

```
ext_if = "re0" # macro for external interface - use tun0 or pppoe0 for PPPoE
int_if = "re1" # macro for internal interface
localnet = $int_if:network
client_out = "{ ssh, domain, pop3, auth, nntp, http,\n              https, cvspserver, 2628, 5999, 8000, 8080 }"
udp_services = "{ domain, ntp }"
icmp_types = "{ echoreq, unreach }"
# ext_if IP address could be dynamic, hence ($ext_if)
match out on $ext_if from $localnet nat-to ($ext_if)
block all
pass quick inet proto { tcp, udp } from $localnet to port $udp_services
pass log inet proto icmp icmp-type $icmp_types
pass inet proto tcp from $localnet port $client_out
```

---

If you are running a PF version equal to OpenBSD 4.6 or earlier, the match rule with `nat-to` instead becomes this:

---

```
nat on $ext_if from $localnet to any -> ($ext_if)
```

---

The only difference that is strictly necessary for your access point to work is to change the definition of `int_if` to match the wireless interface. In our example, this means the line should now read as follows:

---

```
int_if = "ral0" # macro for internal interface
```

---

More than likely, you will also want to set up `dhcpcd` to serve addresses and other relevant network information to clients after they have associated with your access point. Setting up `dhcpcd` is fairly straightforward if you read the man pages.

That's all there is to it. This configuration gives you a functional BSD access point, with at least token security (actually more like a Keep Out! sign) via WEP encryption, or a slightly more robust link-level encryption with WPA. If you need to support FTP, copy the `ftp-proxy` configuration from the machine you set up in Chapter 3 and make changes similar to those you made for the rest of the rule set.

### ***Access Points with Three or More Interfaces***

If your network design dictates that your access point is also the gateway for a wired local network, or even several wireless networks, you need to make some minor changes to your rule set. Instead of just changing the value of the `int_if` macro, you might want to add another (descriptive) definition for the wireless interface, such as the following:

---

```
air_if = "ral0"
```

---

Your wireless interfaces are likely to be on separate subnets, so it might be useful to have a separate rule for each of them to handle the NAT configuration. Here's an example for OpenBSD 4.7 and newer systems:

---

```
match out on $ext_if from $air_if:network nat-to ($ext_if)
```

---

And here's one on pre-OpenBSD 4.7 PF versions:

---

```
nat on $ext_if from $air_if:network to any -> ($ext_if) static-port
```

---

Depending on your policy, you might also want to adjust your `localnet` definition, or at least include `$air_if` in your `pass` rules where appropriate. And once again, if you need to support FTP, a separate `pass` with redirection for the wireless network to `ftp-proxy` may be in order.

### ***Handling IPsec, VPN Solutions***

You can set up virtual private networks (VPNs), using built-in IPsec tools, OpenSSH, or other tools. However, with the relatively poor security profile of wireless networks in general, you are likely to want to set up some additional security.

The options fall roughly into three categories:

**SSH** If your VPN is based on SSH tunnels, the baseline rule set already contains all the filtering you need. Your tunneled traffic will be indistinguishable from other SSH traffic to the packet filter.

**IPsec with UDP key exchange (IKE/ISAKMP)** Several IPsec variants depend critically on key exchange via proto udp port 500 and use proto udp port 4500 for NAT Traversal (NAT-T). You need to let this traffic through in order to let the flows become established. Almost all implementations also depend critically on letting ESP protocol traffic (protocol number 50) pass between the hosts with the following:

---

```
pass proto esp from $source to $target
```

---

**Filtering on IPsec encapsulation interfaces** With a properly configured IPsec setup, you can set up PF to filter on the encapsulation interface `enc0` itself with the following:<sup>4</sup>

---

```
pass on enc0 proto ipencap from $source to $target keep state (if-bound)
```

---

See Appendix A for references to some of the more useful literature on the subject.

## **The Client Side**

As long as you have BSD clients, setup is extremely easy. The steps involved in connecting a BSD machine to a wireless network are quite similar to the ones we just went through to set up a wireless access point. On OpenBSD, the configuration centers on the `hostname.if` file for the wireless interface. On FreeBSD, the configuration centers on `rc.conf`, but will most likely involve a few other files, depending on your exact configuration.

### **OpenBSD Setup**

Starting with the OpenBSD case, in order to connect to the WEP access point we just configured, your OpenBSD clients need a `hostname.if` (for example, `/etc/hostname.ral0`) configuration file with these lines:

---

```
up media autoselect mode 11g chan 1 nwid unwiredbsd nwkey 0x1deadbeef9
dhcpc
```

---

The first line sets the link-level parameters in more detail than usually required. Only `up` and the `nwid` and `nwkey` parameters are strictly necessary. In almost all cases, the driver will associate with the access point on the appropriate channel and in the best available mode. The second line calls for a DHCP configuration, and in practice, causes the system to run a `dhclient` command to retrieve TCP/IP configuration information.

4. In OpenBSD 4.8 the encapsulation interface became a cloneable interface, and you can configure several separate `enc` interfaces. All `enc` interfaces become members of the `enc` interface group.

If you chose to go with the WPA configuration, the file will look like this instead:

---

```
up media autoselect mode 11g chan 1 nwid unwiredbsd wpa wpapsk `wpa-psk unwiredbsd 0x1deadbeef9`  
dhcp
```

---

Again, the first line sets the link-level parameters, where the crucial ones are the network selection and encryption parameters `nwid`, `wpa`, and `wpapsk`. You can try omitting the `mode` and `chan` parameters; in almost all cases, the driver will associate with the access point on the appropriate channel and in the best available mode.

If you want to try out the configuration commands from the command line before committing the configuration to your `/etc/hostname.if` file, the command to set up a client for the WEP network is as follows:

---

```
$ sudo ifconfig ral0 up mode 11b chan 1 nwid unwiredbsd nwkey 0x1deadbeef9
```

---

The `ifconfig` command should complete without any output. You can then use `ifconfig` to check that the interface was successfully configured. The output should look something like this:

---

```
$ ifconfig ral0  
ral0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500  
      lladdr 00:25:9c:72:cf:60  
      priority: 4  
      groups: wlan  
      media: IEEE802.11 autoselect (OFDM54 mode 11g)  
      status: active  
      ieee80211: nwid unwiredbsd chan 1 bssid 00:25:9c:72:cf:60 nwkey <not displayed> 100dBm  
      inet6 fe80::225:9cff:fe72:cf60%ral0 prefixlen 64 scopeid 0x2
```

---

Note that the `ieee80211` line displays the network name and channel, along with a few other parameters. The information displayed here should match what you entered on the `ifconfig` command line.

Here is the command to configure your OpenBSD client to connect to the WPA network:

---

```
$ sudo ifconfig ral0 nwid unwiredbsd wpa wpapsk `wpa-psk unwiredbsd 0x1deadbeef9`
```

---

The command should complete without any output. If you use `ifconfig` again to check the interface status, the output will look something like this:

---

```
$ ifconfig ral0  
ral0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500  
      lladdr 00:25:9c:72:cf:60  
      priority: 4  
      groups: wlan  
      media: IEEE802.11 autoselect (OFDM54 mode 11g)  
      status: active
```

---

```
ieee80211: nwid unwiredbsd chan 1 bssid 00:25:9c:72:cf:60 wpapsk <not displayed>
wpaprotos wpa1,wpa2 wpaakms psk wpaciphers tkip,ccmp wpagroupcipher tkip 100dBm
inet6 fe80::225:9cff:fe72:cf60%ral0 prefixlen 64 scopeid 0x2
```

---

Check that the ieee80211: line displays the correct network name and sensible WPA parameters.

Once you are satisfied that the interface is configured at the link level, use the dhclient command to configure the interface for TCP/IP, like this:

---

```
$ sudo dhclient ral0
```

---

The dhclient command should print a summary of its dialogue with the DHCP server that looks something like this:

```
DHCPREQUEST on ral0 to 255.255.255.255 port 67
DHCPREQUEST on ral0 to 255.255.255.255 port 67
DHCPACK from 10.50.90.1 (00:25:9c:72:cf:60)
bound to 10.50.90.11 -- renewal in 1800 seconds.
```

---

## FreeBSD Setup

On FreeBSD, you may need to do a bit more work than is necessary with OpenBSD. Depending on your kernel configuration, you may need to add the relevant module load lines to */boot/loader.conf*. On one of my test systems, */boot/loader.conf* looks like this:

---

```
if_rum_load="YES"
wlan_scan_ap_load="YES"
wlan_scan_sta_load="YES"
wlan_wep_load="YES"
wlan_ccmp_load="YES"
wlan_tkip_load="YES"
```

---

With the relevant modules loaded, you can join the WEP network we configured earlier by issuing the following command:

---

```
$ sudo ifconfig wlan create wlan0 rum0 ssid unwiredbsd wepmode on wepkey 0x1deadbeef9 up
```

---

Then issue this command:

---

```
$ sudo dhclient wlan0
```

---

For a more permanent configuration, create a *start\_if.rum0* file (replace *rum0* with the name of the physical interface if it differs) with content like this:

---

```
wlans_rum0="wlan0"
create_args_wlan0="wlandev rum0 ssid unwiredbsd wepmode on wepkey 0x1deadbeef9 up"
ifconfig_wlan0="DHCP"
```

---

If you want to join the WPA network, you need to set up `wpa_supplicant` and change your network interface settings slightly. For the WPA access point, connect with the following configuration in your `start_if.rum0` file:

---

```
wlans_rum0="wlan0"
create_args_wlan0="wlandev rum0"
ifconfig_wlan0="WPA"
```

---

You also need an `/etc/wpa_supplicant.conf` file that contains the following:

---

```
network={
    ssid="unwiredbsd"
    psk="0x1deadbeef9"
}
```

---

Finally, add a second `ifconfig_wlan0` line in `rc.conf`, to ensure that `dhclient` runs correctly.

---

```
ifconfig_wlan0="DHCP"
```

---

Other WPA networks may require additional options. After a successful configuration, the `ifconfig` output should display something like this:

---

```
rum0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 2290
      ether 00:24:1d:9a:bf:67
      media: IEEE 802.11 Wireless Ethernet autoselect mode 11g
      status: associated
wlan0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> metric 0 mtu 1500
      ether 00:24:1d:9a:bf:67
      inet 10.50.90.16 netmask 0xffffffff broadcast 10.50.90.255
      media: IEEE 802.11 Wireless Ethernet OFDM/36Mbps mode 11g
      status: associated
      ssid unwiredbsd channel 1 (2412 Mhz 11g) bssid 00:25:9c:72:cf:60
      country US authmode WPA2/802.11i privacy ON deftxkey UNDEF
      TKIP 2:128-bit txpower 0 bmiss 7 scanvalid 450 bgscan bgscanintvl 300
      bgscanidle 250 roam:rssi 7 roam:rate 5 protmode CTS roaming MANUAL
```

---

## Guarding Your Wireless Network with `authpf`

Security professionals tend to agree that even though WEP encryption offers little protection, it's just barely enough to signal to would-be attackers that you do not intend to let all and sundry use your network resources. Using WPA increases security somewhat, at the cost of some complexity.

The configurations we built so far in this chapter are functional. Both the WEP and WPA configurations will let all reasonably configured wireless clients connect, and that may be a problem in itself, since that configuration does not have any real support built in for letting you decide who uses your network.

As mentioned earlier, MAC address filtering is not really a solid defense against attackers because it's just too easy to change a MAC address. The OpenBSD developers chose a radically different approach to this problem when they introduced `authpf` in OpenBSD version 3.1. Instead of tying access to a hardware identifier such as the network card's MAC address, they decided that the robust and highly flexible user authentication mechanisms already in place were more appropriate for the task. The user shell `authpf` lets the system load PF rules on a per-user basis, effectively deciding which user gets to do what.

To use `authpf`, you create users with the `authpf` program as their shell. In order to get network access, the user logs in to the gateway using SSH. Once the user successfully completes SSH authentication, `authpf` loads the rules you have defined for the user or the relevant class of users.

These rules, which apply to the IP address the user logged in from, stay loaded and in force for as long as the user stays logged in via the SSH connection. Once the SSH session is terminated, the rules are unloaded, and in most scenarios, all non-SSH traffic from the user's IP address is denied. With a reasonable setup, only traffic originated by authenticated users will be let through.

**NOTE** *On OpenBSD, `authpf` is one of the login classes offered by default, as you will notice the next time you create a user with `adduser`.*

For systems where the `authpf` login class is not available by default, you may need to add the following lines to your `/etc/login.conf` file:

---

```
authpf:\n    :welcome=/etc/motd.authpf:\n    :shell=/usr/sbin/authpf:\n    :tc=default:
```

---

The next couple of sections contain a few examples that may or may not fit your situation directly, but that I hope will give you ideas you can use.

## A Basic Authenticating Gateway

Setting up an authenticating gateway with `authpf` involves creating and maintaining a few files besides your basic `pf.conf`. The main addition is `authpf.rules`. The other files are fairly static entities that you will not be spending much time on once they have been created.

Start by creating an empty `/etc/authpf/authpf.conf` file. This file needs to be there in order for `authpf` to work, but it doesn't actually need any content, so creating an empty file with `touch` is appropriate.

The other relevant bits of `/etc/pf.conf` follow. First, here are the interface macros:

---

```
ext_if = "re0"\nint_if = "ath0"
```

---

In addition, if you define a table called `<authpf_users>`, authpf will add the IP addresses of authenticated users to the table:

---

```
table <authpf_users> persist
```

---

If you need to run NAT, the rules that take care of the translation could just as easily go in `authpf.rules`, but keeping them in the `pf.conf` file does not hurt in a simple setup like this:

---

```
pass out on $ext_if from $localnet nat-to ($ext_if)
```

---

Here's pre-OpenBSD 4.7 syntax:

---

```
nat on $ext_if from $localnet to any -> ($ext_if)
```

---

Next, we create the authpf anchor, where rules from `authpf.rules` are loaded once the user authenticates:

---

```
anchor "authpf/*"
```

---

For pre-OpenBSD 4.7 authpf versions, several anchors were required, so the corresponding section would be as follows:

---

```
nat-anchor "authpf/*"  
rdr-anchor "authpf/*"  
binat-anchor "authpf/*"  
anchor "authpf/*"
```

---

This brings us to the end of the required parts of a `pf.conf` file for an authpf setup.

For the filtering part, we start with the `block all default`, and then add the pass rules we need. The only essential item at this point is to let SSH traffic pass on the internal network:

---

```
pass quick on $int_if inet proto { tcp, udp } to $int_if port ssh
```

---

From here on out, it really is up to you. Do you want to let your clients have name resolution before they authenticate? If so, put the pass rules for the TCP and UDP service `domain` in your `pf.conf` file, too.

For a relatively simple and egalitarian setup, you could include the rest of our baseline rule set, changing the pass rules to allow traffic from the addresses in the `<authpf_users>` table, rather than any address in your local network:

---

```
pass quick inet proto { tcp, udp } from <authpf_users> to port $udp_services  
pass inet proto tcp from <authpf_users> to port $client_out
```

---

For a more differentiated setup, you could put the rest of your rule set in `/etc/authpf/authpf.rules` or per-user rules in customized `authpf.rules` files in each user's directory under `/etc/authpf/users/`. If your users generally need some protection, your general `/etc/authpf/authpf.rules` could have content like this:

---

```
client_out = "{ ssh, domain, pop3, auth, nntp, http, https }"
udp_services = "{ domain, ntp }"
pass quick inet proto { tcp, udp } from $user_ip to port $udp_services
pass inet proto tcp from $user_ip to port $client_out
```

---

The macro `user_ip` is built into `authpf` and expands to the IP address from which the user authenticated. These rules will apply to any user who completes authentication at your gateway.

A nice and relatively easy addition to implement is special-case rules for users with different requirements than your general user population. If an `authpf.rules` file exists in the user's directory under `/etc/authpf/users/`, the rules in that file will be loaded for the user. This means that your naïve user Peter who only needs to surf the Web and have access to a service that runs on a high port on a specific machine could get what he needs with a `/etc/authpf/users/peter/authpf.rules` file like this:

---

```
client_out = "{ domain, http, https }"
pass inet from $user_ip to 192.168.103.84 port 9000
pass quick inet proto { tcp, udp } from $user_ip to port $client_out
```

---

On the other hand, Peter's colleague Christina runs OpenBSD and generally knows what she is doing, even if she sometimes generates traffic to and from odd ports. You could give her free rein by putting this in `/etc/authpf/users/christina/authpf.rules`:

---

```
pass from $user_ip os = "OpenBSD" to any
```

---

This means Christina can do pretty much anything she likes over TCP as long as she authenticates from her OpenBSD machines.

## ***Wide Open but Actually Shut***

In some settings, it makes sense to set up your network to be open and unencrypted at the link level, while enforcing some restrictions via `authpf`. The next example is very similar to Wi-Fi zones you may encounter in airports or other public spaces, where anyone can associate to the access points and get an IP address, but any attempt at accessing the Web will be redirected to one specific web page until the user has cleared some sort of authentication.<sup>5</sup>

This `pf.conf` file is again built on our baseline, with two important additions to the basic `authpf` setup: a macro and a redirection.

---

5. Thanks to Vegard Engen for the idea and showing me his configuration, which is preserved here in spirit, if not all details.

---

```
ext_if = "re0"
int_if = "ath0"
auth_web="192.168.27.20"
dhcp_services = "{ bootps, bootpc }" # DHCP server + client
table <authpf_users> persist
pass in quick on $int_if proto tcp from ! <authpf_users> to port http rdr-to $auth_web
match out on $ext_if from $int_if:network nat-to ($ext_if)
anchor "authpf/*"
block all
pass quick on $int_if inet proto { tcp, udp } to $int_if port $dhcp_services
pass quick inet proto { tcp, udp } from $int_if:network to any port domain
pass quick on $int_if inet proto { tcp, udp } to $int_if port ssh
```

---

For older authpf versions, use this file instead:

---

```
ext_if = "re0"
int_if = "ath0"
auth_web="192.168.27.20"
dhcp_services = "{ bootps, bootpc }" # DHCP server + client
table <authpf_users> persist
rdr pass on $int_if proto tcp from ! <authpf_users> to any port http -> $auth_web
nat on $ext_if from $localnet to any -> ($ext_if)
nat-anchor "authpf/*"
rdr-anchor "authpf/*"
binat-anchor "authpf/*"
anchor "authpf/*"
block all
pass quick on $int_if inet proto { tcp, udp } to $int_if port $dhcp_services
pass quick inet proto { tcp, udp } from $int_if:network to port domain
pass quick on $int_if inet proto { tcp, udp } to $int_if port ssh
```

---

The `auth_web` macro and the redirection make sure all web traffic from addresses that are not in the `<authpf_users>` table leads all nonauthenticated users to a specific address. At that address, you set up a web server that serves up whatever you need. This could range from a single page with instructions on who to contact in order to get access to the network all the way up to a system that accepts credit cards and handles user creation.

Note that in this setup, name resolution will work, but all surfing attempts will end up at the `auth_web` address. Once the users clear authentication, you can add general rules or user-specific ones to the `authpf.rules` files as appropriate for your situation.

# 5

## **BIGGER OR TRICKIER NETWORKS**



In this chapter, we will build on the material in previous chapters to meet the real-life challenges of both large and small networks with relatively demanding applications or users.

The sample configurations in this chapter are based on the assumption that your packet-filtering setups will need to accommodate services you run on your local network. We will mainly look at this challenge from a Unix perspective, focusing on SSH, email, and web services (with some pointers on how to take care of other services).

This chapter is about the things to do when you need to combine packet filtering with services that must be accessible outside your local network. How much this complicates your rule sets will depend on your network design, and to a certain extent, on the number of routable addresses you have available. We'll begin with configurations for official, routable addresses, and then move on to situations with as few as one routable address and the PF-based work-arounds that make the services usable even under these restrictions.

## A Web Server and Mail Server on the Inside— Routable Addresses

How complicated is your network? How complicated does it need to be?

We'll start with the baseline scenario of the sample clients from Chapter 3 that we set up behind a basic PF firewall, with access to a range of services hosted elsewhere but no services running on the local network. These clients get three new neighbors: a mail server, a web server, and a file server. In this scenario, we use official, routable addresses, since it makes life a little easier. Another advantage of this approach is that with routable addresses, we can let two of the new machines run DNS for our *example.com* domain: one as the master and the other as an authoritative slave.<sup>1</sup>

**NOTE** *For DNS, it always makes sense to have at least one authoritative slave server somewhere outside your own network (in fact, some top-level domains will not let you register a domain without it). You may also want to arrange for a backup mail server to be hosted elsewhere. Keep these things in mind as you build your network.*

At this stage, we keep the physical network layout fairly simple. We put the new servers into the same local network as the clients, possibly in a separate server room, but certainly on the same network segment or switch as the clients. Conceptually, the new network looks something like Figure 5-1.

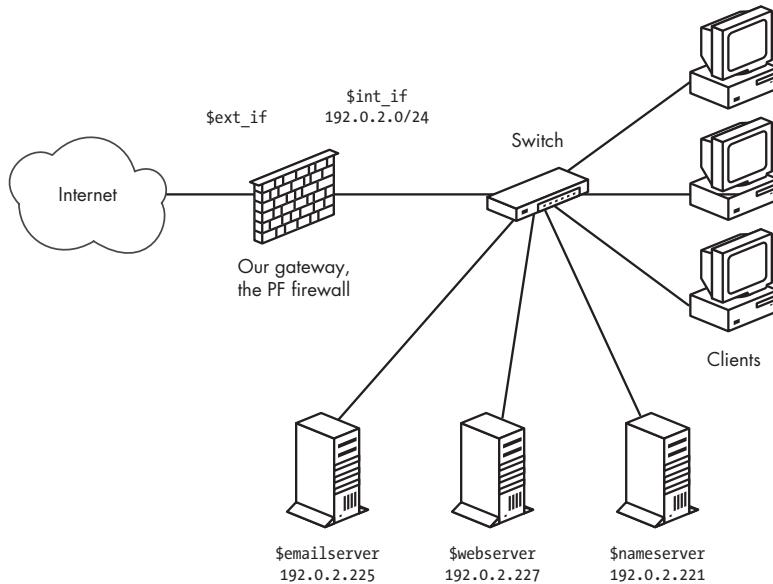


Figure 5-1: A basic network with servers and clients on the inside

1. In fact, the *example.com* network here lives in the 192.0.2.0/24 block, which is set aside in RFC 3330 as reserved for example and documentation use. We use this address range mainly to differentiate from the NAT examples elsewhere in this book, which use addresses in the “private” RFC 1918 address space.

With the basic parameters for the network in place, we can start setting up a sensible rule set for handling the services we need. Once again, we start from the baseline rule set and add a few macros for readability.

The macros we need come rather naturally from the specifications:

- Web server:

---

```
webserver = "192.0.2.227"
```

---

- Web server services:

---

```
webports = "{ http, https }"
```

---

- Mail server:

---

```
emailserver = "192.0.2.225"
```

---

- Mail server services:

---

```
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
```

---

- Name servers:

---

```
nameservers = "{ 192.0.2.221, 192.0.2.223 }"
```

---

We assume that the file server does not need to be accessible to the outside world, unless we choose to set it up with a service that needs to be visible outside the local network, such as an authoritative slave name server for our domain. Then, with the macros in hand, we add the `pass` rules. Starting with the web server, we make it accessible to the world with the following:

---

```
pass proto tcp to $webserver port $webports
```

---

On a similar note, we let the world talk to the mail server:

---

```
pass proto tcp to $emailserver port $email
```

---

This lets clients anywhere have the same access as the ones in your local network, including a few mail-retrieval protocols that run without encryption. That's common enough in the real world, but you might want to consider your options if you are setting up a new network.

For the mail server to be useful, it needs to be able to send mail to hosts outside the local network, too:

---

```
pass log proto tcp from $emailserver to port smtp
```

---

Keeping in mind that the rule set starts with a `block all` rule, this means that only the mail server is allowed to initiate SMTP traffic from the local network to the rest of the world. If any of the other hosts on the network need to send email to the outside world or receive email, they need to use the

designated mail server. This could be a good way to ensure, for example, that you make it as hard as possible for any spam-sending zombie machines that might turn up in your network to actually deliver their payloads.

## IS SYNPROXY WORTH THE TROUBLE?

Over the years, the synproxy state option has received a lot of attention as a possible bulwark against ill-intentioned traffic from the outside. Specifically, the synproxy state option was intended to protect against SYN-flood attacks that could lead to resource exhaustion at the back end.

It works like this: When a new connection is created, PF normally lets the communication partners handle the connection setup themselves, simply passing the packets on if they match a pass rule. With synproxy enabled, PF handles the initial connection setup and hands over the connection to the communication partners only once it is properly established, essentially creating a buffer between the communication partners. The SYN proxying is slightly more expensive than the default keep state, but not necessarily noticeably so on reasonably scaled equipment.

The potential downsides become apparent in load-balancing setups where a SYN-proxying PF could accept connections that the back end is not ready to accept, in some cases short-circuiting the redundancy by setting up connections to other hosts than what the load-balancing logic would have selected. This is especially apparent in protocols like SMTP, where the built-in redundancy dictates that if a primary mail exchanger is not accepting connections, you should try a secondary instead.

When considering a setup where synproxy seems attractive, keep these issues in mind and analyze the potential impact on your setup that would come from adding synproxy to the mix. If you conclude that SYN proxying is needed, simply tack on synproxy state at the end of the rules that need the option.

Finally, the name servers need to be accessible to clients outside our network who look up the information about *example.com* and any other domains for which we answer authoritatively:

---

```
pass inet proto { tcp, udp } to $nameservers port domain
```

---

Having integrated all the services that need to be accessible from the outside world, our rule set ends up looking roughly like this:

---

```
ext_if = "ep0" # macro for external interface - use tun0 or pppoe0 for PPPoE
int_if = "ep1" # macro for internal interface
localnet = $int_if:network
webserver = "192.0.2.227"
webports = "{ http, https }"
emailserver = "192.0.2.225"
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
nameservers = "{ 192.0.2.221, 192.0.2.223 }"
client_out = "{ ssh, domain, pop3, auth, nntp, http,\n              https, cvspserver, 2628, 5999, 8000, 8080 }"
udp_services = "{ domain, ntp }"
icmp_types = "{ echoreq, unreach }"
block all
pass quick inet proto { tcp, udp } from $localnet to port $udp_services
```

```
pass log inet proto icmp all icmp-type $icmp_types
pass inet proto tcp from $localnet to port $client_out
pass inet proto { tcp, udp } to $nameservers port domain
pass proto tcp to $webserver port $webports
pass log proto tcp to $emailserver port $email
pass log proto tcp from $emailserver to port smtp
```

---

This is still a fairly simple setup, but unfortunately, it has one potentially troubling security disadvantage. The way this network is designed, the servers that offer services to the world at large are all *in the same local network* as your clients, and you would need to restrict any internal services to only local access. In principle, this means that an attacker would need to compromise only one host in your local network to gain access to any resource there, putting the miscreant on equal footing with any user in your local network. Depending on how well each machine and resource are protected from unauthorized access, this could be anything from a minor annoyance to a major headache.

In the next section, we will look at some options for segregating the services that need to interact with the world at large from the local network.

## **A Degree of Separation: Introducing the DMZ**

In the previous section, you saw how to set up services on your local network and make them selectively available to the outside world through a sensible PF rule set. For more fine-grained control over access to your internal network, as well as the services you need to make visible to the rest of the world, add a degree of physical separation. Even a separate virtual local area network (VLAN) will do nicely.

Achieving the physical and logical separation is fairly easy: Simply move the machines that run the public services to a separate network, attached to a separate interface on the gateway. The net effect is a separate network that is not quite part of your local network, but not entirely in the public part of the Internet either. Conceptually, the segregated network looks like Figure 5-2.

**NOTE** *Think of this little network as a zone of relative calm between the territories of hostile factions. It is no great surprise that a few years back, someone coined the phrase demilitarized zone, or DMZ for short, to describe this type of configuration.*

For address allocation, you can segment off an appropriately sized chunk of your official address space for the new DMZ network. Alternatively, you can move those parts of your network that do not have a specific need to run with publicly accessible and routable addresses into a NAT environment. Either way, you end up with at least one more interface on which to filter. As you will see later, it is possible to run a DMZ setup in all-NAT environments as well, if you are really short of official addresses.

The adjustments to the rule set itself do not need to be extensive. If necessary, you can change the configuration for each interface. The basic rule set logic remains, but you may need to adjust the definitions of the macros (`webserver`, `mailserver`, `nameservers`, and possibly others) to reflect your new network layout.

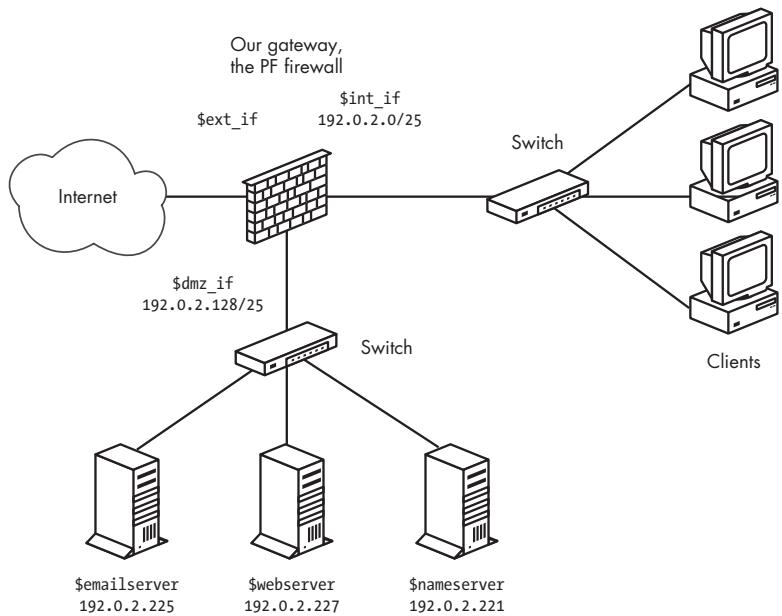


Figure 5-2: A network with the servers in a DMZ

In our example, we could choose to segment off the part of our address range where we have already placed our servers. If we leave some room for growth, we can set up the new `$dmz_if` on a /25 subnet with a network address and netmask of `192.0.2.128/255.255.255.128`. This leaves us with the range `192.0.2.129` through `192.0.2.254` as the usable address range for hosts in the DMZ. With that configuration and no changes in the IP addresses assigned to the servers, you do not really need to touch the rule set at all for the packet filtering to work after setting up a physically segregated DMZ. That is a nice side effect, which could be due to either laziness or excellent long-range planning. Either way, it underlines the importance of having a sensible address-allocation policy in place.

It might be useful to tighten up your rule set by editing your `pass` rules so the traffic to and from your servers is allowed to pass only on the interfaces that are actually relevant to the services:

---

```

pass in on $ext_if proto { tcp, udp } to $nameservers port domain
pass in on $int_if proto { tcp, udp } from $localnet to $nameservers \
    port domain
pass out on $dmz_if proto { tcp, udp } to $nameservers port domain
pass in on $ext_if proto tcp to $webserver port $webports
pass in on $int_if proto tcp from $localnet to $webserver port $webports
pass out on $dmz_if proto tcp to $webserver port $webports
pass in log on $ext_if proto tcp to $mailserver port smtp
pass in log on $int_if proto tcp from $localnet to $mailserver port $email
pass out log on $dmz_if proto tcp to $mailserver port smtp
pass in on $dmz_if from $mailserver to port smtp
pass out log on $ext_if proto tcp from $mailserver to port smtp

```

---

You could choose to make the other pass rules that reference your local network interface-specific, too, but if you leave them intact, they will continue to work.

## **Sharing the Load: Redirecting to a Pool of Addresses**

Once you have set up services to be accessible to the world at large, one likely scenario is that over time, one or more of your services will grow more sophisticated and resource-hungry, or simply attract more traffic than you feel comfortable serving from a single server.

There are a number of ways to make several machines share the load of running a service, including ways to fine-tune the service itself. For the network-level load-balancing, PF offers the basic functionality you need via redirection to tables or address pools. In fact, you can implement a form of load balancing without even touching your pass rules.

Take the web server in our example. We already have the macro for the public IP address (`webserver = "192.0.2.227"`), which, in turn, is associated with the hostname that your users have bookmarked, possibly `www.example.com`. When the time comes to share the load, set up the required number of identical, or at least equivalent, servers, and then alter your rule set slightly to introduce the redirection. First, define a table that holds the addresses for your web server pool:

---

```
table <webpool> persist { 192.0.2.214, 192.0.2.215, 192.0.2.216, 192.0.2.217 }
```

---

Then perform the redirection:

---

```
match in on $ext_if proto tcp to $webserver port $webports \
    rdr-to <webpool> round-robin
```

---

Unlike the redirections in earlier examples, such as the FTP proxy in Chapter 3, this rule sets up all members of the `webpool` table as potential redirection targets for incoming connections intended for the `webports` ports on the `webserver` address. Each incoming connection that matches this rule is redirected to one of the addresses in the table, spreading the load across several hosts. You may choose to retire the original web server once the switch to this redirection is complete, or let it be absorbed in the new web server pool.

On PF versions earlier than OpenBSD 4.7, the equivalent rule is as follows:

---

```
rdr on $ext_if proto tcp to $webserver port $webports -> <webpool> round-robin
```

---

In both cases, the `round-robin` option means that PF shares the load between the machines in the pool by cycling through the table of redirection address sequentially.

When it is essential for accesses from each individual source address to always go to the same host in the back end (for example, if the service depends on client- or session-specific parameters that will be lost if new connections hit a different host in the back end), you can add the `sticky-address` option

to make sure that new connections from a client are always redirected to the same machine behind the redirection as the initial connection. The downside to this option is that PF needs to maintain source-tracking data for each client, and the default value for maximum source nodes tracked is set at 10,000, which may be a limiting factor. (See Chapter 9 for advice on adjusting this and similar limit values.)

When even load distribution is not an absolute requirement, selecting the redirection address at `random` may be appropriate:

---

```
match in on $ext_if proto tcp to $webserver port $webports \
    rdr-to <webpool> random
```

---

**NOTE** *On pre-OpenBSD 4.7 PF versions, the `random` option is not supported for redirection to tables or lists of addresses.*

Even organizations with large pools of official, routable addresses have opted to introduce NAT between their load-balanced server pools and the Internet at large. This technique works equally well in various NAT-based set-ups, but moving to NAT offers some additional possibilities and challenges.

### ***Getting Load Balancing Right with relayd***

After you have been running a while with load balancing via round-robin redirection, you may notice that the redirection does not automatically adapt to external conditions. For example, unless special steps are taken, if a host in the list of redirection targets goes down, traffic will still be redirected to the IP addresses in the list of possibilities.

Clearly, a monitoring solution is needed. Fortunately, the OpenBSD base system provides one. The relay daemon `relayd`<sup>2</sup> interacts with your PF configuration, providing the ability to weed out nonfunctioning hosts from your pool. Introducing `relayd` into your setup, however, may require some minor changes to your rule set.

The `relayd` daemon works in terms of two main classes of services that it refers to as *redirects* and *relays*. It expects to be able to add or subtract hosts' IP addresses to or from the PF tables it controls. The daemon interacts with your rule set through a special-purpose anchor named `relayd` (and in pre-OpenBSD 4.7 versions, also a redirection anchor, `rdr-anchor`, with the same name).

To see how we can make our sample configuration work a little better by using `relayd`, we'll look back at the load-balancing rule set. Starting from the top of your `pf.conf` file, add the anchor for `relayd` to fill in:

---

```
anchor "relayd/*"
```

---

2. Originally introduced in OpenBSD 4.1 under the name `hoststated`, the daemon has seen active development (mainly by Reyk Floeter and Pierre-Yves Ritschard) over several years, including a few important changes to the configuration syntax, and was renamed `relayd` in time for the OpenBSD 4.3 release.

On pre-OpenBSD 4.7 versions, you also need the redirection anchor, like this:

---

```
rdr-anchor "relayd/*"  
anchor "relayd/*"
```

---

In the load-balancing rule set, we had the following definition for our web server pool:

---

```
table webpool persist { 192.0.2.214, 192.0.2.215, 192.0.2.216, 192.0.2.217 }
```

---

It has this `match` rule to set up the redirection:

---

```
match in on $ext_if proto tcp to $webserver port $webports \  
rdr-to <webpool> round-robin
```

---

Or on pre-OpenBSD 4.7 versions, you would use the following:

---

```
rdr on $ext_if proto tcp to $webserver port $webports -> <webpool> round-robin
```

---

To make this configuration work slightly better, we remove the redirection and the table, and let `relayd` handle the redirection by setting up its own versions inside the anchor. (Do not remove the `pass` rule, however, because your rule set will still need to have a `pass` rule that lets traffic flow to the IP addresses in `relayd`'s tables.)

Once the `pf.conf` parts have been taken care of, we turn to `relayd`'s own `relayd.conf` configuration file. The syntax in this configuration file is similar enough to `pf.conf` to make it fairly easy to read and understand. First, we add the macro definitions we will be using later:

---

```
web1="192.0.2.214"  
web2="192.0.2.215"  
web3="192.0.2.216"  
web4="192.0.2.217"  
webserver="192.0.2.227"  
sorry_server="192.0.2.200"
```

---

All of these correspond to definitions we could have put in a `pf.conf` file. The default checking interval in `relayd` is 10 seconds, which means that a host could be down for almost 10 seconds before it is taken offline. Being cautious, we'll set the checking interval to 5 seconds to minimize visible downtime, with the following line:

---

```
interval 5 # check hosts every 5 seconds
```

---

Now we make a table called `webpool` that uses most of the macros:

---

```
table <webpool> { $web1, $web2, $web3, $web4 }
```

---

For reasons we will return to shortly, we define one other table:

---

```
table <sorry> { $sorry_server }
```

---

At this point, we are ready to set up the redirect:

---

```
redirect www {  
    listen on $webserver port 80 sticky-address  
    tag relayd  
    forward to <webpool> check http "/status.html" code 200 timeout 300  
    forward to <sorry> timeout 300 check icmp  
}
```

---

This says that connections to port 80 should be redirected to the members of the `webpool` table. The `sticky-address` option has the same effect here as with the `rdr-to` in PF rules: New connections from the same source IP address (within the time interval defined by the `timeout` value) are redirected to the same host in the back-end pool as the previous ones.

The `relayd` daemon should check to see if a host is available by asking it for the file `/status.html`, using the protocol HTTP, and expecting the return code to be equal to 200. This is the expected result for a client asking a running web server for a file it has available.

No big surprises so far, right? The `relayd` daemon will take care of excluding hosts from the table if they go down. But what if all the hosts in the `webpool` table go down? Fortunately, the developers thought of that too, and introduced the concept of backup tables for services. This is the last part of the definition for the `www` service, with the table `sorry` as the backup table: The hosts in the `sorry` table take over if the `webpool` table becomes empty. This means that you need to configure a service that is able to offer a “Sorry, we’re down” message in case all the hosts in your web pool fail.

With all of the elements of a valid `relayd` configuration in place, you can enable your new configuration. Reload your PF rule set, and then start `relayd`. If you want to check your configuration before actually starting `relayd`, you can use the `-n` command-line option to `relayd`:

---

```
$ sudo relayd -n
```

---

If your configuration is correct, `relayd` displays the message `configuration OK` and exits.

To actually start the daemon, we do not need any command-line flags, so the following sequence reloads your edited PF configuration and enables `relayd`.

---

```
$ sudo pfctl -f /etc/pf.conf  
$ sudo relayd
```

---

With a correct configuration, both commands will silently start, without displaying any messages. You can check that relayd is running with top or ps. In both cases, you will find three relayd processes, roughly like this:

---

```
$ ps aux | grep relayd
relayd  9153  0.0  0.1  776 1424 ??  S      7:28PM   0:00.01 relayd: pf update engine
(relayd)
relayd  6144  0.0  0.1  776 1440 ??  S      7:28PM   0:00.02 relayd: host check engine
(relayd)
root    3217  0.0  0.1  776 1416 ??  Is     7:28PM   0:00.01 relayd: parent (relayd)
```

---

In almost all cases, you will want to enable relayd at startup. You do that by adding this line to your *rc.conf.local* file:

---

```
relayd_flags="" # for normal use: ""
```

---

However, once the configuration is enabled, most of your interaction with relayd will happen through the relayctl administration program. In addition to letting you monitor status, relayctl lets you reload the relayd configuration and selectively disable or enable hosts, tables, and services. You can even view service status interactively, like this:

---

```
$ sudo relayctl show summary
Id      Type          Name           Avlblty Status
1       redirect      www            active
1       table         webpool:80    active (2 hosts)
1       host          192.0.2.214  100.00% up
2       host          192.0.2.215  0.00%  down
3       host          192.0.2.216  100.00% up
4       host          192.0.2.217  0.00%  down
2       table         sorry:80    active (1 hosts)
5       host          127.0.0.1    100.00% up
```

---

In this example, the web pool is seriously degraded, with only two of four hosts up and running. Fortunately, the backup table is still functioning. All tables are active with at least one host up. For tables that no longer have any members, the Status column changes to empty. Asking relayctl for host information shows the status information in a host-centered format:

---

```
$ sudo relayctl show hosts
Id      Type          Name           Avlblty Status
1       table         webpool:80    active (3 hosts)
1       host          192.0.2.214  100.00% up
                           total: 11340/11340 checks
2       host          192.0.2.215  0.00%  down
                           total: 0/11340 checks, error: tcp connect failed
3       host          192.0.2.216  100.00% up
                           total: 11340/11340 checks
4       host          192.0.2.217  0.00%  down
                           total: 0/11340 checks, error: tcp connect failed
```

---

---

2	table	sorry:80	active (1 hosts)
5	host	127.0.0.1	100.00% up
		total: 11340/11340 checks	

---

If you need to take a host out of the pool for maintenance (or any time-consuming operation), you can use `relayctl` to disable it like this:

---

```
$ sudo relayctl host disable 192.0.2.217
```

---

In most cases, the operation will display `command succeeded` to indicate that the operation completed successfully. Once you have completed maintenance and put the machine online, you can reenable it as part of `relayd`'s pool with this command:

---

```
$ sudo relayctl host enable 192.0.2.217
```

---

Again, you should see the message `command succeeded` almost immediately to indicate that the operation was successful.

In addition to the basic load-balancing demonstrated here, `relayd` has been extended in recent OpenBSD versions to offer several features that make it attractive in more complex settings. For example, it can now handle Layer 7 proxying or relaying functions for HTTP and HTTPS, including protocol handling with header append and rewrite, URL path append and rewrite, and even session and cookie handling. The protocol handling needs to be tailored to your application. For example, the following is a simple HTTPS relay for load-balancing the encrypted web traffic from clients to the web servers.

---

```
http protocol "httpssl" {
    header append "$REMOTE_ADDR" to "X-Forwarded-For"
    header append "$SERVER_ADDR:$SERVER_PORT" to "X-Forwarded-By"
    header change "Keep-Alive" to "$TIMEOUT"
    query hash "sessid"
    cookie hash "sessid"
    path filter "*command=*" from "/cgi-bin/index.cgi"

    ssl { sslv2, ciphers "MEDIUM:HIGH" }
    tcp { nodelay, sack, socket buffer 65536, backlog 128 }
}
```

---

This protocol handler definition demonstrates a range of simple operations on the HTTP headers, and sets both SSL parameters and specific TCP parameters to optimize connection handling. The header options operate on the protocol headers, inserting the values of the variables by either appending to existing headers (`append`) or changing the content to a new value (`change`).

The URL and cookie hashes are used by the load-balancer to select to which host in the target pool the request is forwarded. The `path filter` specifies that any get request, including the first quoted string as a substring of the second, is to be dropped. The `ssl` options specify that only SSL version 2

ciphers are accepted, with key lengths in the medium-to-high range; in other words, 128 bits or more.<sup>3</sup> Finally, the `tcp` options specify `nodelay` to minimize delays, specify the use of the selective acknowledgment method (RFC 2018), and set the socket buffer size and the maximum allowed number of pending connections the load-balancer keeps track of. These options are examples only; in most cases, your application will perform well with these settings at their default values.

The relay definition using the protocol handler follows a pattern that should be familiar from the earlier definition of the `www` service:

---

```
relay wwwssl {
    # Run as a SSL accelerator
    listen on $webserver port 443 ssl
    protocol "httpssl"
    table <webhosts> loadbalance check ssl
}
```

---

Still, your SSL-enabled web applications will likely benefit from a slightly different set of parameters.

**NOTE** *We've added a `check ssl`, assuming that each member of the `webhosts` table is properly configured to complete an SSL handshake. Depending on your application, it may be useful to look into keeping all SSL processing in `relayd`, thus offloading the encryption-handling tasks from the back ends.*

Finally, for CARP-based failover of the hosts running `relayd` on your network (see “Redundancy and Failover: CARP and `pfsync`” on page 119), `relayd` can be configured to support CARP interaction by setting the CARP demotion counter for the specified interface groups at shutdown or startup.

Like all parts of the OpenBSD system, `relayd` comes with informative man pages. For the angles and options not covered here (there are a few), dive into the man pages for `relayd`, `relayd.conf`, and `relayctl`, and start experimenting to find just the configuration you need.

## A Web Server and Mail Server on the Inside—the NAT Version

Let's backtrack a little and begin again with the baseline scenario where the sample clients from Chapter 3 get three new neighbors: a mail server, a web server, and a file server. This time around, externally visible addresses are either not available or too expensive, and running several other services on a machine that is primarily a firewall is not desirable. This means we are back to the situation where we do our NAT at the gateway. Fortunately, the redirection mechanisms in PF make it relatively easy to keep servers on the inside of a gateway that performs NAT.

The network specifications are the same as for the `example.com` setup we just worked through: We need to run a web server that serves up data in clear-text (`http`) and encrypted (`https`), and we want a mail server that sends and

---

3. See the OpenSSL man page for further explanation of cipher-related options.

receives email while letting clients inside and outside the local network use a number of well-known submission and retrieval protocols. In short, we want pretty much the same features as in the setup from the previous section, but with only one routable address.

Of the three servers, only the web server and the mail server need to be visible to the outside world, so we add macros for their IP addresses and services to the Chapter 3 rule set:

---

```
webserver = "192.168.2.7"
webports = "{ http, https }"
emailserver = "192.168.2.5"
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
```

---

With only one routable address and the servers hidden in NATed address space, we need to set up rules at the gateway that redirect the traffic we want our servers to handle. We could define a set of `match` rules to set up the redirection, and then address the `block` or `pass` question in a separate set of rules later, like this:

---

```
match in on $ext_if proto tcp to $ext_if port $webports rdr-to $webserver
match in on $ext_if proto tcp to $ext_if port $email rdr-to $emailserver

pass proto tcp to $webserver port $webports
pass proto tcp to $emailserver port $email
pass proto tcp from $emailserver to port smtp
```

---

This combination of `match` and `pass` rules is very close to the way things were done in pre-OpenBSD 4.7 PF versions, and if you are upgrading from a previous version, this is the kind of quick edit that could bridge the syntax gap quickly. But you could also opt to go for the new style, and write this slightly more compact version instead:

---

```
pass in on $ext_if inet proto tcp to $ext_if port $webports rdr-to $webserver
pass in on $ext_if inet proto tcp to $ext_if port $email rdr-to $mailserver
pass on $int_if inet proto tcp to $webserver port $webports
pass on $int_if inet proto tcp to $mailserver port $email
```

---

Note the use of `pass` rules with `rdr-to`. This combination of filtering and redirection will help make things easier in a little while, so try this combination for now.

On pre-OpenBSD 4.7 PF, the rule set will be quite similar, except in the way that we handle the redirections.

---

```
webserver = "192.168.2.7"
webports = "{ http, https }"
emailserver = "192.168.2.5"
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
```

---

```
rdr on $ext_if proto tcp to $ext_if port $webports -> $webserver  
rdr on $ext_if proto tcp to $ext_if port $email -> $emailserver  
  
pass proto tcp to $webserver port $webports  
pass proto tcp to $emailserver port $email  
pass proto tcp from $emailserver to any port smtp
```

---

## DMZ with NAT

With an all-NAT setup, the pool of available addresses to allocate for a DMZ is likely to be larger than in our previous example, but the same principles apply. When you move the servers off to a physically separate network, you will need to check that your rule set's macro definitions are sane and adjust the values if necessary.

Just as in the routable addresses case, it might be useful to tighten up your rule set by editing your `pass` rules so the traffic to and from your servers is allowed to pass on only the interfaces that are actually relevant to the services:

```
pass in on $ext_if inet proto tcp to $ext_if port $webports rdr-to $webserver  
pass in on $int_if inet proto tcp from $localnet to $webserver port $webports  
pass out on $dmz_if proto tcp to $webserver port $webports  
pass in log on $ext_if inet proto tcp to $ext_if port $email \  
    rdr-to $mailserver  
pass in log on $int_if proto tcp from $localnet to $mailserver port $email  
pass out log on $dmz_if proto tcp to $mailserver port smtp  
pass in on $dmz_if from $mailserver to port smtp  
pass out log on $ext_if proto tcp from $mailserver to port smtp
```

---

The version for pre-OpenBSD 4.7 PF differs in some details, with the redirection still in separate rules:

```
pass in on $ext_if proto tcp to $webserver port $webports  
pass in on $int_if proto tcp from $localnet to $webserver port $webports  
pass out on $dmz_if proto tcp to $webserver port $webports  
pass in log on $ext_if proto tcp to $mailserver port smtp  
pass in log on $int_if proto tcp from $localnet to $mailserver port $email  
pass out log on $dmz_if proto tcp to $mailserver port smtp  
pass in on $dmz_if from $mailserver to port smtp  
pass out log on $ext_if proto tcp from $mailserver to port smtp
```

---

You could create specific `pass` rules that reference your local network interface, but if you leave the existing `pass` rules intact, they will continue to work.

## Redirection for Load Balancing

The redirection-based load-balancing rules from the previous example work equally well in a NAT regime, where the public address is the gateway's external interface and the redirection addresses are in a private range.

Here's the webpool definition:

---

```
table <webpool> persist { 192.168.2.7, 192.168.2.8, 192.168.2.9, 192.168.2.10 }
```

---

The main difference between the routable-address case and the NAT version is that after you have added the webpool definition, you edit the existing pass rule with redirection, which then becomes this:

---

```
pass in on $ext_if inet proto tcp to $ext_if port $webports rdr-to <webpool> round-robin
```

---

Or for pre-OpenBSD 4.7 PF versions, use this:

---

```
rdr on $ext_if proto tcp to $ext_if port $webports -> <webpool> round-robin
```

---

From that point on, your NATed DMZ behaves much like the one with official, routable addresses.

### ***Back to the Single NATed Network***

It may surprise you to hear that there are cases where setting up a small network is more difficult than working with a large one. For example, returning to the situation where the servers are on the same physical network as the clients, the basic NATed configuration works very well—up to a point. In fact, everything works brilliantly as long as all you are interested in is getting traffic from hosts outside your local network to reach your servers.

Here is the full configuration:

---

```
ext_if = "re0" # macro for external interface - use tun0 or pppoe0 for PPPoE
int_if = "re1" # macro for internal interface
localnet = $int_if:network
# for ftp-proxy
proxy = "127.0.0.1"
icmp_types = "{ echoreq, unreach }"
client_out = "{ ssh, domain, pop3, auth, nntp, http, https, \
               446, cvspserver, 2628, 5999, 8000, 8080 }"
udp_services = "{ domain, ntp }"
webserver = "192.168.2.7"
webports = "{ http, https }"
emailserver = "192.168.2.5"
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
# NAT: ext_if IP address could be dynamic, hence ($ext_if)
match out on $ext_if from $localnet nat-to ($ext_if)
block all
# for ftp-proxy: Remember to put the following line, uncommented, in your
# /etc/rc.conf.local to enable ftp-proxy:
# ftpproxy_flags=""
anchor "ftp-proxy/*"
pass in quick proto tcp to port ftp rdr-to $proxy port 8021
pass out proto tcp from $proxy to port ftp
pass quick inet proto { tcp, udp } to port $udp_services
pass proto tcp to port $client_out
# allow out the default range for traceroute(8):
```

```
# "base+nhops*nqueries-1" (33434+64*3-1)
pass out on $ext_if inet proto udp to port 33433 >< 33626 keep state
# make sure icmp passes unfettered
pass inet proto icmp icmp-type $icmp_types from $localnet
pass inet proto icmp icmp-type $icmp_types to $ext_if
pass in on $ext_if inet proto tcp to $ext_if port $webports rdr-to $webserver
pass in on $ext_if inet proto tcp to $ext_if port $email rdr-to $mailserver
pass on $int_if inet proto tcp to $webserver port $webports
pass on $int_if inet proto tcp to $mailserver port $email
```

---

The last four rules here are the ones that interest us the most. If you try to reach the services on the official address from hosts in your own network, you will soon see that the requests for the redirected services from machines in your local network most likely never reach the external interface. This is because all the redirection and translation happens on the external interface. The gateway receives the packets from your local network on the internal interface, with the destination address set to the external interface's address. The gateway recognizes the address as one of its own and tries to handle the request as if it were directed at a local service; as a consequence, the redirections do not quite work from the inside.

The equivalent part to those last four lines of the preceding rule set for pre-OpenBSD 4.7 systems looks like this:

```
rdr on $ext_if proto tcp to $ext_if port $webports -> $webserver
rdr on $ext_if proto tcp to $ext_if port $email -> $mailserver

pass proto tcp to $webserver port $webports
pass proto tcp to $mailserver port $email
pass proto tcp from $mailserver to any port smtp
```

---

Fortunately, several work-arounds for this particular problem are possible. The problem is common enough that the *PF User Guide* lists four different solutions to the problem,<sup>4</sup> including moving your servers to a DMZ as described earlier. Since this is a PF book, we will concentrate on a PF-based solution (actually a pretty terrible work-around), which consists of treating the local network as a special case for our redirection and NAT rules.

We need to intercept the network packets originating in the local network and handle those connections correctly, making sure that any return traffic is directed to the communication partner who actually originated the connection. This means that in order for the redirections to work as expected from the local network, we need to add special-case redirection rules that mirror the ones designed to handle requests from the outside. First, here are the pass rules with redirections for OpenBSD 4.7 and newer:

```
pass in on $ext_if inet proto tcp to $ext_if port $webports rdr-to $webserver
pass in on $ext_if inet proto tcp to $ext_if port $email rdr-to $mailserver
pass in log on $int_if inet proto tcp from $int_if:network to $ext_if \
    port $webports rdr-to $webserver
```

---

4. See the “Redirection and Reflection” section in the *PF User Guide* (<http://www.openbsd.org/faq/pf/rdr.html#reflect>).

```
pass in log on $int_if inet proto tcp from $int_if:network to $ext_if \
    port $email rdr-to $mailserver
match out log on $int_if proto tcp from $int_if:network to $webserver \
    port $webports nat-to $int_if
pass on $int_if inet proto tcp to $webserver port $webports
match out log on $int_if proto tcp from $int_if:network to $mailserver \
    port $email nat-to $int_if
pass on $int_if inet proto tcp to $mailserver port $email
```

---

The first two rules are identical to the original ones. The next two intercept the traffic from the local network and the `rdr-to` actions in both rewrite the destination address, much as the corresponding rules do for the traffic that originates elsewhere. The `pass on $int_if` rules serve the same purpose as in the earlier version.

The `match` rules with `nat-to` are there as a routing work-around. Without them, the `webserver` and `mailserver` hosts would route return traffic for the redirected connections directly back to the hosts in the local network, where the traffic would not match any outgoing connection. With the `nat-to` in place, the servers consider the gateway as the source of the traffic, and will direct return traffic back the same path it came originally. The gateway matches the return traffic to the states created by connections from the clients in the local network, and applies the appropriate actions to return the traffic to the correct clients.

The equivalent rules for pre-OpenBSD 4.7 versions are at first sight a bit more confusing, but the end result is the same.

---

```
rdr on $int_if proto tcp from $localnet to $ext_if port $webports -> $webserver
rdr on $int_if proto tcp from $localnet to $ext_if port $email -> $emailserver
no nat on $int_if proto tcp from $int_if to $localnet
nat on $int_if proto tcp from $localnet to $webserver port $webports -> $int_if
nat on $int_if proto tcp from $localnet to $emailserver port $email -> $int_if
```

---

This way, we twist the redirections and the address translation logic to do what we need, and we do not need to touch the `pass` rules at all. (I've had the good fortune to witness via email or IRC the reactions of several network admins at the point when the truth about this five-line reconfiguration sank in.)

## Filtering on Interface Groups

Your network could have several subnets that may never need to interact with your local network except for some common services like email, web, file, and print. How you handle the traffic from and to such subnets depends on how your network is designed. One useful approach is to treat each less-privileged network as a separate local network attached to its own separate interface on a common filtering gateway, and then give it a rule set that allows only the desired direct interaction with the neighboring networks attached to the main gateway.

You can make your PF configuration more manageable and readable by grouping logically similar interfaces into *interface groups* and apply filtering rules to the groups rather than the individual interfaces. Interface groups, as implemented via the `ifconfig group` option, originally appeared in OpenBSD 3.6 and have been adopted in FreeBSD 7.0 onward.

All configured network interfaces can be configured to belong to one or more groups. Some interfaces automatically belong to one of the default groups. For example, all IEEE 802.11 wireless network interfaces belong to the `wlan` group, while interfaces associated with the default routes belong to the `egress` group. Fortunately, an interface can be a member of several groups, and you can add interfaces to interface groups via the appropriate `ifconfig` command, as in this example:

---

```
# ifconfig sis2 group untrusted
```

---

The equivalent under OpenBSD would be in the `hostname.sis2` file or the `ifconfig_sis2=` line in the `rc.conf` file on FreeBSD 7.0 or later.

Where it makes sense, you can then treat the interface group much the same as you would handle a single interface in filtering rules:

---

```
pass in on untrusted to any port $webports  
pass out on egress to any port $webports
```

---

If by now you're thinking that in most, if not all, the rule set examples up to this point it would be possible to filter on the group `egress` instead of the macro `$ext_if`, you have grasped an important point. It could be a useful exercise to go through any existing rule sets you have and see what using interface groups can do to help readability even further. Remember that an interface group can have one or more members.

Note that filtering on interface groups makes it possible to write essentially hardware-independent rule sets. As long as your `hostname.if` files or `ifconfig_if=` lines put the interfaces in the correct groups, rule sets that consistently filter on interface groups will be fully portable between machines that may or may not have identical hardware configurations.

On systems where the interface group feature is not available, you may be able to achieve some of the same effects via creative use of macros, like this:

---

```
untrusted = "{ ath0 ath1 wlo ep0 }"  
egress = "sk0"
```

---

## The Power of Tags

In some networks, the decision of where a packet should be allowed to pass cannot be made to map easily to criteria like subnet and service. The fine-grained control the site's policy demands could make the rule set complicated and potentially hard to maintain.

Fortunately, PF offers yet another mechanism for classification and filtering in the form of packet *tagging*. The useful way to implement packet tagging is to tag incoming packets that match a specific pass rule, and then let the packets pass elsewhere based on which identifiers the packet is tagged with. In OpenBSD 4.6 and later, it is even possible to have separate `match` rules that tag according to the match criteria, leaving decisions on passing, redirecting, or taking other actions to rules later in the rule set.

One example could be the wireless access points we set up in Chapter 4, which we could reasonably expect to inject traffic into the local network with an apparent source address equal to the access point's `$ext_if` address. In that scenario, a useful addition to the rule set of a gateway with several of these access points might be the following (assuming, of course, that definitions of the `wifi_allowed` and `wifi_ports` macros fit the site's requirements).

---

```
wifi = "{ 10.0.0.115, 10.0.0.125, 10.0.0.135, 10.0.0.145 }"
pass in on $int_if from $wifi to $wifi_allowed port $wifi_ports tag wifigood
pass out on $ext_if tagged wifigood
```

---

As the complexity of rule set grows, consider using `tag` in incoming `match` and `pass` rules to make your rule set readable and easier to maintain.

Tags are sticky, and once a packet has been tagged by a matching rule, the tag stays, which means that a packet can have a tag even if it was not applied by the last matching rule. However, a packet can have only one tag at any time. If a packet matches several rules that apply tags, the tag will be overwritten with a new one by each new matching tag rule.

For example, you could set several tags on incoming traffic via a set of `match` or `pass` rules, supplemented by a set of `pass` rules that determine where packets pass out based on the tags set on the incoming traffic.

## The Bridging Firewall

An Ethernet *bridge* consists of two or more interfaces configured to forward Ethernet frames transparently, and which are not directly visible to the upper layers, such as the TCP/IP stack. In a filtering context, the bridge configuration is often considered attractive because it means that the filtering can be performed on a machine that does not have its own IP addresses. If the machine in question runs OpenBSD or a similarly capable operating system, it can still filter and redirect traffic.

The main advantage of such a setup is that attacking the firewall itself is more difficult.<sup>5</sup> The disadvantage is that all admin tasks must be performed at the firewall's console, unless you configure a network interface that is reachable via a secured network of some kind, or even a serial console. It also follows that bridges with no IP address configured cannot be set as the gateway

---

5. How much security this actually adds is a matter of occasional heated debate on mailing lists such as `openbsd-misc` and other networking-oriented lists. Reading up on the pros and cons as perceived by core OpenBSD developers can be entertaining as well as enlightening.

for a network and cannot run any services on the bridged interfaces. Rather, you can think of a bridge as an intelligent bulge on the network cable, which can filter and redirect.

A few general caveats apply to using firewalls implemented as bridges:

- The interfaces are placed in promiscuous mode, which means that they will receive (and to some extent process) every packet on the network.
- Bridges operate on the Ethernet level and, by default, forward all types of packets, not just TCP/IP.
- The lack of IP addresses on the interfaces makes some of the more effective redundancy features, such as CARP, unavailable.

The method for configuring bridges differs in some details among operating systems. The following examples are very basic and do not cover all possible wrinkles, but should be enough to get you started.

### ***Basic Bridge Setup on OpenBSD***

The OpenBSD GENERIC kernel contains all the necessary code to configure bridges and filter on them. Unless you have compiled a custom kernel without the bridge code, the setup is quite straightforward.

**NOTE**

*On OpenBSD 4.7 and newer, the brconfig command no longer exists. All bridge configuration and related functionality was merged into ifconfig for the OpenBSD 4.7 release.*

To set up a bridge with two interfaces on the command line, you first create the bridge device. The first device of a kind is conventionally given the sequence number 0, so we create the `bridge0` device with the following command:

---

```
$ sudo ifconfig bridge0 create
```

---

Before the next `ifconfig` command (brconfig on OpenBSD 4.6 and earlier), use `ifconfig` to check that the prospective member interfaces (in our case, `ep0` and `ep1`) are up, but not assigned IP addresses. Next, configure the bridge by entering the following:

---

```
$ sudo ifconfig bridge0 add ep0 add ep1 blocknonip ep0 blocknonip ep1 up
```

---

On OpenBSD 4.6 and earlier, enter this:

---

```
$ sudo brconfig bridge0 add ep0 add ep1 blocknonip ep0 blocknonip ep1 up
```

---

The OpenBSD `ifconfig` command (brconfig on OpenBSD 4.6 and earlier) contains a fair bit of filtering code itself. In this example, we use the `blocknonip` option for each interface to block all non-IP traffic.

**NOTE** The OpenBSD ifconfig command offers its own set of filtering options in addition to other configuration options. The bridge(4) and ifconfig(8) man pages provide further information. Since it operates on the Ethernet level, it is possible to use ifconfig to specify filtering rules that let the bridge filter on MAC addresses. Using these filtering capabilities, it is also possible to let the bridge tag packets for further processing in your PF rule set via the tagged keyword. For tagging purposes, a bridge with one member interface will do. This functionality and the keywords were also present in brconfig(4) on OpenBSD 4.6 and earlier.

To make the configuration permanent, create or edit /etc/hostname.ep0 and enter the following line:

---

up

---

For the other interface, /etc/hostname.ep1 should contain the same line:

---

up

---

Finally, enter the bridge setup in /etc/hostname.bridge0 (/etc/bridgename.bridge0 on OpenBSD 4.6 and earlier):

---

add ep0 add ep1 blocknonip ep0 blocknonip ep1 up

---

Your bridge should now be up, and you can go on to create the PF filter rules.

### **Basic Bridge Setup on FreeBSD**

For FreeBSD, the procedure is a little more involved than on OpenBSD. In order to be able to use bridging, your running kernel must include the `if_bridge` module. The default kernel configurations build this module, so under ordinary circumstances, you can go directly to creating the interface. To compile the bridge device into the kernel, add the following line in the kernel configuration file:

---

device if\_bridge

---

You can also load the device at boot time by putting the following line in the `/etc/loader.conf` file.

---

if\_bridge\_load="YES"

---

Create the bridge by entering this:

---

\$ sudo ifconfig bridge0 create

---

Creating the `bridge0` interface also creates a set of bridge-related `sysctl` values:

---

```
$ sudo sysctl net.link.bridge  
net.link.bridge.ipfw: 0  
net.link.bridge.pfil_member: 1  
net.link.bridge.pfil_bridge: 1  
net.link.bridge.ipfw_arp: 0  
net.link.bridge.pfil_onlyip: 1
```

---

It is worth checking that these `sysctl` values are available. If they are, it is confirmation that the bridge has been enabled. If they are not, go back and see what went wrong and why.

**NOTE** *These values apply to filtering on the bridge interface itself. You do not need to touch them, since IP-level filtering on the member interfaces (the ends of the pipe) is enabled by default.*

Before the next `ifconfig` command, check that the prospective member interfaces (in our case, `ep0` and `ep1`) are up but have not been assigned IP addresses, and then configure the bridge by entering this:

---

```
$ sudo ifconfig bridge0 addm ep0 addm ep1 up
```

---

To make the configuration permanent, add the following lines to `/etc/rc.conf`:

---

```
ifconfig_ep0="up"  
ifconfig_ep1="up"  
cloned_interfaces="bridge0"  
ifconfig_bridge0="addm ep0 addm ep1 up"
```

---

This means your bridge is up, and you can go on to creating the PF filter rules. See the `if_bridge(4)` man page for further FreeBSD-specific bridge information.

### **Basic Bridge Setup on NetBSD**

On NetBSD, the default kernel configuration does not have the filtering bridge support compiled in. You need to compile a custom kernel with the following option added to the kernel configuration file. Once you have the new kernel with the bridge code in place, the setup is straightforward.

---

options	BRIDGE_IPF	# bridge uses IP/IPv6 pfil hooks too
---------	------------	--------------------------------------

---

To create a bridge with two interfaces on the command line, first create the `bridge0` device:

---

```
$ sudo ifconfig bridge0 create
```

---

Before the next `brconfig` command, use `ifconfig` to check that the prospective member interfaces (in our case, `ep0` and `ep1`) are up but have not been assigned IP addresses. Then configure the bridge by entering this:

---

```
$ sudo brconfig bridge0 add ep0 add ep1 up
```

---

Next, enable the filtering on the `bridge0` device:

---

```
$ sudo brconfig bridge0 ipf
```

---

To make the configuration permanent, create or edit `/etc/ifconfig.ep0` and enter the following line.

---

```
up
```

---

For the other interface, `/etc/ifconfig.ep1` should contain the same line:

---

```
up
```

---

Finally, enter the bridge setup in `/etc/ifconfig.bridge0`:

---

```
create  
!add ep0 add ep1 up
```

---

Your bridge should now be up, and you can go on to create the PF filter rules. For further information, see the PF on NetBSD documentation at <http://www.netbsd.org/Documentation/network/pf.html>.

## The Bridge Rule Set

Here is the `pf.conf` file for a bulge-in-the-wire version of the baseline rule set we started with in this chapter. The network changes slightly, as shown in Figure 5-3.

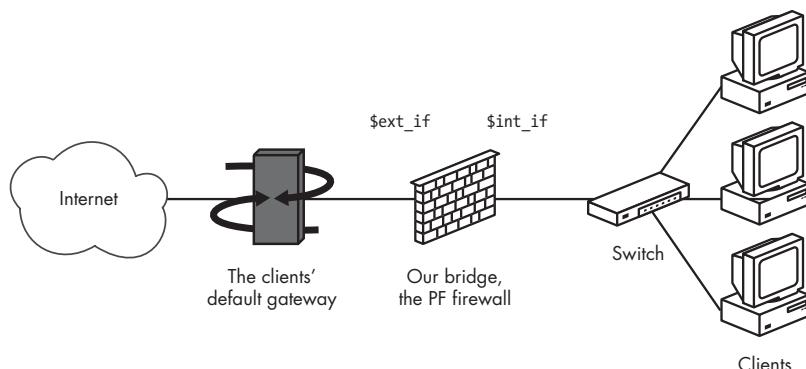


Figure 5-3: A network with a bridge firewall

The machines in the local network share a common default gateway, which is not the bridge, but could be placed either inside or outside the bridge.

---

```
ext_if = ep0
int_if = ep1
localnet= "192.0.2.0/24"
webserver = "192.0.2.227"
webports = "{ http, https }"
emailserver = "192.0.2.225"
email = "{ smtp, pop3, imap, imap3, imaps, pop3s }"
nameservers = "{ 192.0.2.221, 192.0.2.223 }"
client_out = "{ ssh, domain, pop3, auth, nntp, http, https, \
                446, cvspserver, 2628, 5999, 8000, 8080 }"
udp_services = "{ domain, ntp }"
icmp_types = "{ echoreq, unreach }"
set skip on $int_if
block all
pass quick on $ext_if inet proto { tcp, udp } from $localnet \
            to port $udp_services
pass log on $ext_if inet proto icmp all icmp-type $icmp_types
pass on $ext_if inet proto tcp from $localnet to port $client_out
pass on $ext_if inet proto { tcp, udp } to $nameservers port domain
pass on $ext_if proto tcp to $webserver port $webports
pass log on $ext_if proto tcp to $emailserver port $email
pass log on $ext_if proto tcp from $emailserver to port smtp
```

---

Significantly more complicated setups are possible. But remember that while redirections will work, you will not be able to run services on any of the interfaces without IP addresses.

## Handling Nonroutable Addresses from Elsewhere

Even with a properly configured gateway to handle filtering and potentially NAT for your own network, you may find yourself in the unenviable position of needing to compensate for other people's misconfigurations.

One depressingly common class of misconfigurations is the kind that lets traffic with nonroutable addresses out to the Internet. Traffic from nonroutable addresses has also played a part in several denial-of-service (DoS) attack techniques, so it's worth considering explicitly blocking traffic from nonroutable addresses from entering your network. One possible solution is outlined here. For good measure, it also blocks any attempt to initiate contact to nonroutable addresses through the gateway's external interface.

---

```
martians = "{ 127.0.0.0/8, 192.168.0.0/16, 172.16.0.0/12, \
              10.0.0.0/8, 169.254.0.0/16, 192.0.2.0/24, \
              0.0.0.0/8, 240.0.0.0/4 }"

block in quick on $ext_if from $martians to any
block out quick on $ext_if from any to $martians
```

---

Here, the `martians` macro denotes the RFC 1918 addresses and a few other ranges mandated by various RFCs not to be in circulation on the open Internet. Traffic to and from such addresses is quietly dropped on the gateway's external interface.

**NOTE** *The `martians` macro could easily be implemented as a table instead, with all of the table advantages as an added bonus for your rule set.*

The specific details of how to implement this kind of protection will vary according to your network configuration and may be part of a wider set of network security measures. Your network design might also dictate that you include or exclude other address ranges than these.

# 6

## TURNING THE TABLES FOR PROACTIVE DEFENSE



In the previous chapter, you saw how you might need to spend considerable time and energy making sure that the services you want to offer will be available even when you have strict packet filtering in place. Now, with your working setup in place, you will soon notice that some services tend to attract a little more unwanted attention than others.

Here's the scenario: You have a network with packet filtering to match your site's needs, including some services that need to be accessible to users from elsewhere. Unfortunately, when services are available, there is a risk that someone will want to exploit them for some sort of mischief.

You will almost certainly have remote login via SSH, as well as SMTP email running on your network—both are tempting targets. In this chapter, we'll look at some ways to make it harder to gain unauthorized access via SSH, and then we'll turn to some of the more effective ways to deny spammers use of your servers.

## Turning Away the Brutes

The Secure Shell service, commonly referred to as SSH, is a fairly crucial service for Unix administrators. It's frequently the main interface to the machine and a favorite target of script kiddie attacks.

### SSH Brute-Force Attacks

If you run an SSH login service that is accessible from the Internet, you've probably seen entries like this in your authentication logs:

---

```
Sep 26 03:12:34 skapet sshd[25771]: Failed password for root from 200.72.41.31 port 40992 ssh2
Sep 26 03:12:34 skapet sshd[5279]: Failed password for root from 200.72.41.31 port 40992 ssh2
Sep 26 03:12:35 skapet sshd[5279]: Received disconnect from 200.72.41.31: 11: Bye Bye
Sep 26 03:12:44 skapet sshd[29635]: Invalid user admin from 200.72.41.31
Sep 26 03:12:44 skapet sshd[24703]: input_userauth_request: invalid user admin
Sep 26 03:12:44 skapet sshd[24703]: Failed password for invalid user admin from 200.72.41.31 port
41484 ssh2
Sep 26 03:12:44 skapet sshd[29635]: Failed password for invalid user admin from 200.72.41.31 port
41484 ssh2
Sep 26 03:12:45 skapet sshd[24703]: Connection closed by 200.72.41.31
Sep 26 03:13:10 skapet sshd[11459]: Failed password for root from 200.72.41.31 port 43344 ssh2
Sep 26 03:13:10 skapet sshd[7635]: Failed password for root from 200.72.41.31 port 43344 ssh2
Sep 26 03:13:10 skapet sshd[11459]: Received disconnect from 200.72.41.31: 11: Bye Bye
Sep 26 03:13:15 skapet sshd[31357]: Invalid user admin from 200.72.41.31
Sep 26 03:13:15 skapet sshd[10543]: input_userauth_request: invalid user admin
Sep 26 03:13:15 skapet sshd[10543]: Failed password for invalid user admin from 200.72.41.31 port
43811 ssh2
Sep 26 03:13:15 skapet sshd[31357]: Failed password for invalid user admin from 200.72.41.31 port
43811 ssh2
Sep 26 03:13:15 skapet sshd[10543]: Received disconnect from 200.72.41.31: 11: Bye Bye
Sep 26 03:13:25 skapet sshd[6526]: Connection closed by 200.72.41.31
```

---

This is what a *brute-force attack* looks like. Someone or something is trying by brute force to find a username and password to use to get into your system.

The simplest response would be to write a *pf.conf* rule that blocks all access, but that leads to another class of problems, including how to let people with legitimate business on your system access it. Setting up your *sshd* to accept only key-based authentication would help, but most likely would not stop the kiddies from trying. You might consider moving the service to another port, but then again, the ones flooding you on port 22 would probably be able to scan their way to port 22222 for a repeat performance.

Since OpenBSD 3.7 (and equivalents), PF has offered a slightly more elegant solution.

### Setting Up an Adaptive Firewall

To thwart brute-force attacks, you can write your *pass* rules so they maintain certain limits on what connecting hosts can do. For good measure, you can banish violators to a table of addresses that you deny some or all access. You

can even choose to drop all existing connections from machines that over-reach your limits. To enable this feature, first set up the table by adding the following line to your configuration before any filtering rules:

---

```
table <bruteforce> persist
```

---

Then, early in your rule set, block brute force, as shown here:

---

```
block quick from <bruteforce>
```

---

Finally, add your pass rule:

---

```
pass inet proto tcp to $localnet port $tcp_services \
    keep state (max-src-conn 100, max-src-conn-rate 15/5, \
        overload <bruteforce> flush global)
```

---

This rule is very similar to what you've seen in earlier examples. The interesting part in this context is the contents of the parentheses, called *state-tracking options*:

- `max-src-conn` is the number of simultaneous connections allowed from one host. In this example, it's set to 100. You may want a slightly higher or lower value, depending on your network's traffic patterns.
- `max-src-conn-rate` is the rate of new connections allowed from any single host. Here, it's set to 15 connections per 5 seconds, denoted as 15/5. Choose a rate that suits your setup.
- `overload <bruteforce>` means that the address of any host that exceeds the preceding limits is added to the table `bruteforce`. Our rule set blocks all traffic from addresses in the `bruteforce` table.

Once a host exceeds any of these limits and is put in the `overload` table, the rule no longer matches traffic from that host. Make sure that overloaders are handled, if only by a default block rule or similar.

- `flush global` says that when a host reaches the limit, all states for its connections are terminated (flushed). The `global` option means that for good measure, `flush` applies to states created by traffic from that host that matches other `pass` rules, too.

As you can imagine, the effect of this tiny addition to the rule set is quite dramatic. After a few tries, brute force ends up in the `bruteforce` table. That means that all their existing connections are terminated (flushed), and any new attempts will be blocked, most likely with `Fatal: timeout before authentication` messages at their end. You have created an *adaptive firewall* that adjusts automatically to conditions in your network and acts on undesirable activity.

**NOTE** These adaptive rules are effective only for protection against the traditional, rapid-fire type of brute-force attempts. The low-intensity, distributed password-guessing attempts that were first identified as such in 2008 and have been recurring ever since (known among other names as *The Hail Mary Cloud*) do not produce traffic that will match these rules.

It's likely that you may want some flexibility in your rule set, and allow a larger number of connections for some services, but would like to be a little more tight-fisted when it comes to SSH. In that case, you could supplement the general-purpose pass rule with something like the following one, early in your rule set:

---

```
pass quick proto { tcp, udp } to port ssh \
    keep state (max-src-conn 15, max-src-conn-rate 5/3, \
        overload <bruteforce> flush global)
```

---

You should be able to find the set of parameters that is just right for your situation by reading the relevant man pages and the *PF User Guide* (see Appendix A).

**NOTE** Remember that these sample rules are intended as illustrations, and your network's needs may be better served by different rules. Setting the number of simultaneous connections or the rate of connections too low may block legitimate traffic with a clear risk of self-inflicted denial of service when the configuration includes many hosts behind a common NATing gateway, and the users on the NATed hosts have legitimate business on the other side of a gateway with strict overload rules.

The state-tracking options and the overload mechanism do not need to apply exclusively to the SSH service, and blocking all traffic from offenders is not always desired. You could, for example, use a rule like this:

---

```
pass proto { tcp, udp } to port $mail_services \
    keep state (max 1500, max-src-conn 100)
```

---

Here, `max` specifies the maximum number of states that can be created for each rule (the number of rules loaded depends on what the `$mail_services` macro expands to) with no `overload` to protect a mail or web service from receiving more connections than it can handle. With this rule, the `max` value determines the maximum number of states that will be created for each resulting rule. Once the limit is reached, new connections will not be allowed until the old ones terminate. Alternatively, you could remove the `max` restriction, add an `overload` part to the rule, and assign offenders to a queue with a minimal bandwidth allocation (see the discussion of ALTQ in Chapter 7 for details on setting up queues).

Some sites use `overload` to implement a multitiered system, where hosts that trip one `overload` rule are transferred to one or more intermediate "probation" tables for special treatment. It can be useful in web contexts to not block traffic from hosts in the `overload` tables outright, but rather redirect all HTTP requests from these hosts to specific web pages (like in the `authpf` example near the end of Chapter 4).

## **Tidying Your Tables with pfctl**

With the overload rules from the previous section in place, you now have an adaptive firewall that automatically detects undesirable behavior and adds offenders' IP addresses to tables. Watching the logs and the tables can be fun in the short run, but since those rules only add to the tables, we run into the next challenge: keeping the content of the tables up to date and relevant.

When you've run a configuration with an adaptive rule set for a while, at some point, you will discover that an IP address one of your overload rules blocked last week due to a brute-force attack was actually a dynamically assigned address, which is now assigned to a different ISP customer with a legitimate reason to communicate with hosts in your network.<sup>1</sup> If your adaptive rules catch a lot of traffic on a busy network, you may also find that the overload tables will grow over time to take up an increasing amount of memory.

The solution is to *expire* table entries—to remove entries that have not been referenced for a certain amount of time. In OpenBSD 4.1, pfctl acquired the ability to expire table entries based on the time since their statistics were last reset.<sup>2</sup> (In almost all instances, this reset time is equal to the time since the table entry was added.) The keyword is `expire`, and the table entry's age is specified in seconds. Here's an example:

---

```
$ sudo pfctl -t bruteforce -T expire 86400
```

---

This command will remove bruteforce table entries that had their statistics reset more than 86,400 seconds (24 hours) ago.

**NOTE** *The choice of 24 hours as the expiry time is a fairly arbitrary one. You should choose a value that you feel is a reasonable amount of time for any problem at the other end to be noticed and fixed. If you have adaptive rules in place, it is a good idea to set up crontab entries to run table expiry at regular intervals with a command much like the preceding one to make sure your tables are kept up to date.*

## **Giving Spammers a Hard Time with spamd**

Email is a fairly essential service that needs special attention due to the large volume of unwanted messages, or *spam*. The volume of unsolicited commercial messages was already a rather painful problem when malware makers discovered that email-borne worms would work and started using email to spread their payload. During the early 2000s, the combined volume of spam and email-borne malware had increased to the point where running an SMTP mail service without some sort of spam countermeasures had become almost unthinkable.

---

1. In a longer-term perspective, it is fairly normal for entire networks and larger ranges of IP addresses to be reassigned to new owners in response to events in the physical, business-oriented world.

2. Before pfctl acquired the ability to expire table entries, table expiry was more likely than not handled by the special-purpose utility `expiretable`. If your pfctl does not have the `expire` option, look for `expiretable` in your package system.

Spam-fighting measures are almost as old as the spam problem itself. The early efforts focused on analysis of the messages' contents (known as *content filtering*) and to some extent on interpretation of the messages' rather trivially forgeable headers such as the purported sender address (`From:`) or the store and forward paths of intermediate deliveries recorded in the `Received:` headers.

When the OpenBSD team designed their spam-fighting solution `spamd`, first introduced with OpenBSD 3.3 in 2003, they instead focused on the network level and the immediate communication partner in the SMTP conversations, combined with any available information about hosts that tried to deliver messages. The developers set out to create a small, simple, and secure program. The early implementation was based almost entirely on creative use of PF tables combined with data from trusted external sources.

**NOTE** *In addition to the OpenBSD spam-deferral daemon, the content-filtering-based anti-spam package SpamAssassin (<http://spamassassin.apache.org/>) also features a program called `spamd`. Both programs are designed to help fight spam, but they take very different approaches to the underlying problem and do not interoperate directly. However, when both programs are correctly configured and running, they complement each other well.*

## **Network-Level Behavior Analysis and Blacklisting**

The original `spamd` design is based on the observation that spammers send a lot of mail, and the likelihood that you are the first person receiving a particular message is incredibly small. In addition, spam is sent via a few spammer-friendly networks and numerous hijacked machines. Both the individual messages and the machines that send them will be reported to blacklist maintainers quickly, and the blacklist data consisting of known spam senders' IP addresses forms the basis for `spamd`'s processing.

When dealing with blacklisted hosts, `spamd` employs a method called *tarpitting*. When the daemon receives an SMTP connection, it presents its banner and immediately switches to a mode where it answers SMTP traffic at the rate of 1 byte per second, using a tiny selection of SMTP commands designed to make sure that mail is never delivered, but rather rejected back into the sender's queue once the message headers have been transferred. The intention is to waste as much time as possible on the sending end, while costing the receiver pretty much nothing. This specific tarpitting implementation with 1-byte SMTP replies is often referred to as *stuttering*. Blacklist-based tarpitting with stuttering was the default mode for `spamd` up to and including OpenBSD 4.0.

**NOTE** *On FreeBSD and NetBSD, `spamd` is not part of the base system, but available through ports and packages as mail/spamd. If you are running PF on FreeBSD 5.n or NetBSD 4.0 or newer, install the port or package now.*

## Setting Up *spamd* in Blacklisting Mode

To set up *spamd* to run in traditional, blacklisting-only mode, you first put a special-purpose table and a matching redirection in *pf.conf*, and then turn your attention to *spamd*'s own *spamd.conf*. *spamd* then hooks into the PF rule set via the table and the redirection.

The following are the *pf.conf* lines for this configuration:

---

```
table <spamd> persist
pass in on $ext_if inet proto tcp from <spamd> to \
{ $ext_if, $localnet } port smtp rdr-to 127.0.0.1 port 8025
```

---

And here is the pre-OpenBSD 4.7 syntax:

---

```
table <spamd> persist
rdr pass on $ext_if inet proto tcp from <spamd> to \
{ $ext_if, $localnet } port smtp -> 127.0.0.1 port 8025
```

---

The table, *<spamd>*, is there to store the IP addresses you import from trusted blacklist sources. The redirection takes care of all SMTP attempts from hosts that are already in the blacklist. *spamd* listens on port 8025 and responds s-l-o-w-l-y (1 byte per second) to all SMTP connections it receives as a result of the redirection. Later on in the rule set, you would have a rule that makes sure legitimate SMTP traffic passes to the mail server.

*spamd.conf* is where you specify the sources of your blacklist data and any exceptions or local overrides you want.

**NOTE** *On OpenBSD 4.0 and earlier (and by extension, ports based on versions prior to OpenBSD 4.1), spamd.conf was in /etc. Beginning with OpenBSD 4.1, spamd.conf is found in /etc/mail instead. The FreeBSD port installs a sample configuration in /usr/local/etc/spamd/spamd.conf.sample.*

Near the beginning of *spamd.conf*, you will notice a line without a # comment sign that looks like *all:\`*. This line specifies the blacklists you will use. Here is an example:

---

```
all:\`
:uatraps:whiteList:
```

---

Add all blacklists that you want to use below the *all:\`* line, separating each with a colon (:). To use whitelists to subtract addresses from your blacklist, add the name of the whitelist immediately after the name of each blacklist, as in *:blacklist:whiteList:*

Next is the blacklist definition:

---

```
uatraps:\`
:blacklist:\`
:msg="SPAM. Your address %A has sent spam within the last 24 hours":\`
:method=http:\`
:file=www.openbsd.org/spamd/traplist.gz
```

---

Following the name (`uatraps`), the first data field specifies the list type—in this case `black`. The `msg` field contains the message to be displayed to black-listed senders during the SMTP dialogue. The `method` field specifies how `spamd-setup` fetches the list data—in this case via HTTP. Other possibilities include fetching via FTP (`ftp`), from a file in a mounted file system (`file`), or via execution of an external program (`exec`). Finally, the `file` field specifies the name of the file `spamd` expects to receive.

The definition of a whitelist follows much the same pattern, but omits the `message` parameter:

---

```
whitelist:\n    :white:\\\n    :method=file:\\\n    :file=/etc/mail/whitelist.txt
```

---

**NOTE** *The suggested blacklists in the default `spamd.conf` could exclude large blocks of the Internet, including several address ranges that claim to cover entire countries. If your site expects to exchange legitimate mail with any of the countries in question, those lists may not be optimal for your setup. Other popular lists have been known to list entire /16 ranges as spam sources, and it is well worth reviewing the details of the list's maintenance policy before putting a blacklist into production.*

Put the lines for `spamd` and the startup parameters you want in your `/etc/rc.conf.local` on OpenBSD, or in `/etc/rc.conf` on FreeBSD or NetBSD. Here is an example:

---

```
spamd_flags="-v -b" # for normal use: "" and see spamd-setup(8)
```

---

Here, we enable `spamd` and set it to run in blacklisting mode with the `-b` flag. In addition, the `-v` flag enables verbose logging, which is useful for keeping track of `spamd`'s activity for debugging purposes.

On FreeBSD, the `/etc/rc.conf` settings that control `spamd`'s behavior are `obspamd_enable`, which should be set to "YES" in order to enable `spamd`, and `obspamd_flags`, where you fill in any command-line options for `spamd`:

---

```
obspamd_enable="YES"\nobspamd_flags="-v -b" # for normal use: "" and see spamd-setup(8)
```

---

**NOTE** *To have `spamd` run in pure blacklist mode on OpenBSD 4.1 or newer, you can achieve the same effect by setting the `spamd_black` variable to "YES", and then restarting `spamd`.*

Once you've finished editing the setup, start `spamd` with the options you want, and complete the configuration with `spamd-setup`. Finally, create a cron job that calls `spamd-setup` to update the blacklist at reasonable intervals. In pure blacklist mode, you can view and manipulate the table contents using `pfctl` table commands.

## **spamd Logging**

By default, `spamd` logs to your general system logs. To send the `spamd` log messages to a separate log file, add an entry like this to `syslog.conf`:

---

<code>!!spamd</code>		
<code>daemon.err;daemon.warn;daemon.info;daemon.debug</code>		<code>/var/log/spamd</code>

---

Once you’re satisfied that `spamd` is running and does what it is supposed to do, you will probably want to add the `spamd` log file to your log rotations, too. After you’ve run `spamd-setup` and the tables are filled, you can view the table contents using `pfctl`.

**NOTE**

*In the sample pf.conf fragment at the beginning of this section, the redirection (rdr-to) rule is also a pass rule. If you opted for a match rule instead (or if you are using an older PF version and chose to write a rdr rule that does not include a pass part), be sure to set up a pass rule to let traffic through to your redirection. You may also need to set up rules to let legitimate email through. However, if you are already running an email service on your network, you can probably go on using your old SMTP pass rules.*

Given a set of reliable and well-maintained blacklists, `spamd` in pure blacklisting mode does a good job of reducing spam. However, with pure blacklisting, you catch traffic only from hosts that have already tried to deliver spam elsewhere, and you put your trust in external data sources to determine which hosts deserve to be tarpitted. For a setup that provides more immediate response to network-level behavior and offers some real gains in spam prevention, consider *greylisting*, which is a crucial part of how the modern `spamd` works.

## ***Greylisting: My Admin Told Me Not to Talk to Strangers***

Greylisting consists mainly of interpreting the current SMTP standards and adding a little white lie to make life easier.

Spammers tend to use other people’s equipment to send their messages, and the software they install without the legal owner’s permission needs to be relatively lightweight in order to run undetected. Unlike legitimate mail senders, spammers typically do not consider any individual message they send to be important. Taken together, this means that typical spam and malware sender software are not set up to interpret SMTP status codes correctly. This is a fact that we can use to our advantage, as Evan Harris proposed in a paper titled “The Next Step in the Spam Control War: Greylisting,” published in 2003.<sup>3</sup>

As Harris noted, when a compromised machine is used to send spam, the sender application tends to try delivery only once, without checking for any results or return codes. Real SMTP implementations interpret SMTP

---

<sup>3</sup>. The original Harris paper and a number of other useful articles and resources can be found at <http://www.greylisting.org/>.

return codes and act on them, and real mail servers will retry if the initial attempt fails with any kind of temporary error.

In his paper, Harris outlined a practical approach:

- On first SMTP contact from a previously unknown communication partner, *do not* receive email on the first delivery attempt, but instead respond with a status code that indicates a temporary local problem, and store the sender IP address for future reference.
- If the sender retries immediately, reply as before with the temporary failure status code.
- If the sender retries after a set minimum amount of time (1 hour, for example) but not more than a maximum waiting period (4 hours, for example), accept the message for delivery and record the sender IP address in your whitelist.

This is the essence of greylisting. And fortunately, you can set up and maintain a greylisting `spamd` on your PF-equipped gateway.

### Setting Up `spamd` in Greylisting Mode

OpenBSD's `spamd` acquired its ability to greylist in OpenBSD 3.5. Beginning with OpenBSD 4.1, `spamd` runs in greylisting mode by default.

In the default greylisting mode, the `spamd` table used for blacklisting, as described in the previous section, becomes superfluous. You can still use blacklists, but `spamd` will use a combination of private data structures for black-list data and the `spamdb` database to store greylisting-related data. A typical set of rules for `spamd` in default mode looks like this:

---

```
table <spamd-white> persist
table <nospamd> persist file "/etc/mail/nospamd"
pass in log on egress proto tcp to port smtp \
    rdr-to 127.0.0.1 port spamd
pass in log on egress proto tcp from <nospamd> to port smtp
pass in log on egress proto tcp from <spamd-white> to port smtp
pass out log on egress proto tcp to port smtp
```

---

This includes the necessary `pass` rules to let legitimate email flow to the intended destinations from your own network. The `<spamd-white>` table is the whitelist, maintained by `spamd`. The hosts in the `<spamd-white>` table have passed the greylisting hurdle, and mail from these machines is allowed to pass to the real mail servers or their content-filtering front ends. In addition, the `nospamd` table is there for you to load addresses of hosts that you do not want to expose to `spamd` processing, and the matching `pass` rule makes sure SMTP traffic from those hosts passes.

In your network, you may want to tighten those rules to pass SMTP traffic only to and from hosts that are allowed to send and receive email via SMTP. We will get back to the `nospamd` table in “Handling Sites That Do Not Play Well with Greylisting” on page 102.

## WHY GREYLISTING WORKS

A significant amount of design and development effort have been put into making essential services such as SMTP email transmission fault-tolerant. In practical terms, this means that the best effort of a service such as SMTP is as close as you can get to having a perfect record for delivering messages. That's why we can rely on greylisting to eventually let us receive email from proper mail servers.

The current standard for Internet email transmission is defined in RFC 5321.\* The following are several excerpts from Section 4.5.4.1, "Sending Strategy":

In a typical system, the program that composes a message has some method for requesting immediate attention for a new piece of outgoing mail, while mail that cannot be transmitted immediately MUST be queued and periodically retried by the sender.

The sender MUST delay retrying a particular destination after one attempt has failed. In general, the retry interval SHOULD be at least 30 minutes; however, more sophisticated and variable strategies will be beneficial when the SMTP client can determine the reason for non-delivery.

Retries continue until the message is transmitted or the sender gives up; the give-up time generally needs to be at least 4–5 days.

Delivering email is a collaborative, best-effort process, and the RFC clearly states that if the site you are trying to send mail to reports that it can't receive at the moment, it is your duty (a *must* requirement) to try again later, giving the receiving server a chance to recover from its problem.

The clever wrinkle to greylisting is that it's a convenient white lie. When we claim to have a temporary local problem, that problem is really the equivalent of "My admin told me not to talk to strangers." Well-behaved senders with valid messages will call again, but spammers won't wait around for the chance to retry, since doing so increases their cost of delivering messages. This is why greylisting still works, and since it's based on strict adherence to accepted standards,<sup>†</sup> false positives are rare.

\* The relevant parts of RFC 5321 are identical to the corresponding parts of RFC 2821, which is obsolete. Some of us were more than a little disappointed that the IETF did not clarify these chunks of the text, now moving forward on the standards track. My reaction (actually it's quite a rant) is at <http://bsdly.blogspot.com/2008/10/ietf-failed-to-account-for-greylisting.html>.

† The relevant RFCs are mainly RFC 1123 and RFC 5321 (which obsoleted the earlier RFC 2821). Remember that temporary rejection is an SMTP fault-tolerance feature.

The following are the equivalent rules in pre-OpenBSD 4.7 syntax:

```
table <spamd-white> persist
table <nospamd> persist file "/etc/mail/nospamd"
rdr pass in log on egress proto tcp to port smtp \
    -> 127.0.0.1 port spamd
pass in log on egress proto tcp from <nospamd> to port smtp
pass in log on egress proto tcp from <spamd-white> to port smtp
pass out log on egress proto tcp to port smtp
```

On FreeBSD, in order to use `spamd` in greylisting mode, you need a file descriptor filesystem (see `man 5 fdescfs`) mounted at `/dev/fd/`. To implement this, add the following line to `/etc/fstab`, and make sure the `fdescfs` code is in your kernel, either compiled in or by loading the module via the appropriate `kldload` command.

---

```
fdescfs /dev/fd fdescfs rw 0 0
```

---

To begin configuring `spamd`, place the lines for `spamd` and the startup parameters you want in `/etc/rc.conf.local`. Here is an example:

---

```
spamd_flags="-v -G 2:4:864" # for normal use: "" and see spamd-setup(8)
```

---

On FreeBSD, the equivalent line should go in `/etc/rc.conf`:

---

```
obspamd_flags="-v -G 2:4:864" # for normal use: "" and see spamd-setup(8)
```

---

You can fine-tune several of the greylisting-related parameters via `spamd` command-line parameters trailing the `-G` option.

This colon-separated list `2:4:864` represents the values `passtime`, `greyexp`, and `whiteexp`:

- `passtime` denotes the minimum number of minutes `spamd` considers a reasonable time before retry. The default is 25 minutes, but here we've reduced it to 2 minutes.
- `greyexp` is the number of hours an entry stays in the greylisted state before it is removed from the database.
- `whiteexp` determines how long a whitelisted entry is kept, in hours. The default values for `greyexp` and `whiteexp` are 4 hours and 864 hours (just over one month), respectively.

## Greylisting in Practice

Users and administrators at sites that implement greylisting tend to agree that greylisting gets rid of most of their spam, with a significant drop in the load on any content filtering they have in place for their mail. We will start by looking at what `spamd`'s greylisting looks like according to log files, and then return with some data.

If you start `spamd` with the `-v` command-line option for verbose logging, your logs will include a few more items of information in addition to IP addresses. With verbose logging, a typical log excerpt looks like this:

---

```
Oct  2 19:53:21 delilah spamd[26905]: 65.210.185.131: connected (1/1), lists: spews1
Oct  2 19:55:04 delilah spamd[26905]: 83.23.213.115: connected (2/1)
Oct  2 19:55:05 delilah spamd[26905]: (GREY) 83.23.213.115: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
Oct  2 19:55:05 delilah spamd[26905]: 83.23.213.115: disconnected after 0 seconds.
Oct  2 19:55:05 delilah spamd[26905]: 83.23.213.115: connected (2/1)
Oct  2 19:55:06 delilah spamd[26905]: (GREY) 83.23.213.115: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
```

---

```
Oct  2 19:55:06 delilah spamd[26905]: 83.23.213.115: disconnected after 1 seconds.
Oct  2 19:57:07 delilah spamd[26905]: (BLACK) 65.210.185.131: <bounce-3C7E40A4B3@branch15.summer-
bargainz.com> -> <adm@dataped.no>
Oct  2 19:58:50 delilah spamd[26905]: 65.210.185.131: From: Auto Insurance Savings
<noreply@branch15.summer-bargainz.com>
Oct  2 19:58:50 delilah spamd[26905]: 65.210.185.131: Subject: Start SAVING MONEY on Auto Insurance
Oct  2 19:58:50 delilah spamd[26905]: 65.210.185.131: To: adm@dataped.no
Oct  2 20:00:05 delilah spamd[26905]: 65.210.185.131: disconnected after 404 seconds. lists: spews1
Oct  2 20:03:48 delilah spamd[26905]: 222.240.6.118: connected (1/0)
Oct  2 20:03:48 delilah spamd[26905]: 222.240.6.118: disconnected after 0 seconds.
Oct  2 20:06:51 delilah spamd[26905]: 24.71.110.10: connected (1/1), lists: spews1
Oct  2 20:07:00 delilah spamd[26905]: 221.196.37.249: connected (2/1)
Oct  2 20:07:00 delilah spamd[26905]: 221.196.37.249: disconnected after 0 seconds.
Oct  2 20:07:12 delilah spamd[26905]: 24.71.110.10: disconnected after 21 seconds. lists: spews1
```

---

The first line is the beginning of a connection from a machine in the spews1 blacklist. The next six lines show the complete records of two connection attempts from another machine, which each time connects as the second active connection. This second machine is not yet in any blacklist, so it is grey-listed. Note the rather curious delivery address (*whitp98zpu.fsf@datadok.no*) in the message that the greylisted machine tries to deliver. There is a useful trick that we'll look at in "Greytrapping" on page 98. The (GREY) and (BLACK) before the addresses indicate greylisting or blacklisting status. Then there is more activity from the blacklisted host, and a little later we see that after 404 seconds (or 6 minutes, 44 seconds), the blacklisted host gives up without completing the delivery.

The remaining lines show a few very short connections, including one from a machine already in a blacklist. This time, though, the machine disconnects too quickly to see any (BLACK) flag at the beginning of the SMTP dialogue, but we see a reference to the list name (spews1) at the end.

Roughly 400 seconds is about the amount of time that naïve, blacklisted spammers hang around (according to data from various sites), and about the time it takes (at the rate of 1 byte per second) to complete the EHLO ... dialogue until *spamd* rejects the message. However, while peeking at the logs, you are likely to find some spammers that hang around significantly longer. For example, in the data from our office gateway, one log entry stood out:

---

```
Dec 11 23:57:24 delilah spamd[32048]: 69.6.40.26: connected (1/1), lists:
spamhaus spews1 spews2
Dec 12 00:30:08 delilah spamd[32048]: 69.6.40.26: disconnected after 1964
seconds. lists: spamhaus spews1 spews2
```

---

This particular machine was already on several blacklists when it made 13 attempts at delivery from December 9 through December 12. The last attempt lasted 32 minutes, 44 seconds, without completing the delivery. Relatively intelligent spam senders drop the connection during the first few seconds, like the ones in the first log fragment. Others give up after around 400 seconds. A few hang on for hours. (The most extreme case we have recorded hung on for 42,673 seconds, which is almost 12 hours.)

## Tracking Your Real Mail Connections: `spamlogd`

Behind the scenes, rarely mentioned and barely documented, is one of `spamd`'s most important helper programs: the `spamlogd` whitelist updater. As the name suggests, `spamlogd` works quietly in the background, logging connections to and from your mail servers to keep your whitelist updated. The idea is to make sure that valid mail sent between hosts you communicate with regularly goes through with a minimum of fuss.

**NOTE** *If you have followed the discussion up to this point, `spamlogd` has probably been started automatically already. However, if your initial `spamd` configuration did not include greylisting, `spamlogd` may not have been started, and you may experience strange symptoms like the greylist and whitelist not being updated properly. Restarting `spamd` after you have enabled greylisting should ensure that `spamlogd` is loaded and available, too.*

In order to perform its job properly, `spamlogd` needs you to log SMTP connections to and from your mail servers, just as we did in the sample rule sets in Chapter 5:

---

```
emailserver = "192.0.2.225"
pass log proto tcp to $emailserver port $email
pass log proto tcp from $emailserver to port smtp
```

---

On OpenBSD 4.1 and higher (and equivalents), you can create several pflog interfaces and specify where rules should be logged. Here's how to separate the data `spamlogd` needs to read from the rest of your PF logs:

1. Create a separate `pflog1` interface using `ifconfig pflog1 create`, or create a `hostname.pflog1` file with just the line up.
2. Change the rules to the following:

---

```
pass log (to pflog1) proto tcp to $emailserver port $email
pass log (to pflog1) proto tcp from $emailserver to port smtp
```

---

3. Add `-1 pflog1` to `spamlogd`'s startup parameters.

This separates the `spamd`-related logging from the rest. (See Chapter 8 for more about logging.)

With the preceding rules in place, `spamlogd` will add the IP addresses that receive email you send to the whitelist. This is not an ironclad guarantee that the reply will pass immediately, but in most configurations, it helps speed things significantly.

## Greytrapping

We know that spam senders rarely use a fully compliant SMTP implementation to send their messages, which is why greylisting works. We also know that spammers rarely check that the addresses they feed to their hijacked machines

are actually deliverable. Combine these facts, and you see that if a greylisted machine tries to send a message to an invalid address in your domain, there is a good chance that the message is spam or malware.

This realization led to the next evolutionary step in `spamd` development—a technique dubbed *greytrapping*. When a greylisted host tries to deliver mail to a known bad address in our domains, the host is added to a locally maintained blacklist called `spamd-greytrap`. Members of the `spamd-greytrap` list are treated to the same 1-byte-per-second tarpitting as members of other blacklists.

Greytrapping as implemented in `spamd` is simple and elegant. The main thing you need as a starting point is `spamd` running in greylisting mode. The other crucial component is a list of addresses in domains your servers handle email for, but only ones that you are sure will never receive legitimate email. The number of addresses in your list is unimportant, but there must be at least one, and the upper limit is mainly defined by how many addresses you wish to add.

Next, you use `spamdb` to feed your list to the greytrapping feature, and sit back and watch. First, a sender tries to send email to an address on your greytrap list and is simply greylisted, as with any sender you have not exchanged email with before. If the same machine tries again, either to the same, invalid address or another address on your greytrapping list, the greytrap is triggered, and the offender is put into `spamd-greytrap` for 24 hours. For the next 24 hours, any SMTP traffic from the greytrapped host will be stuttered, with 1-byte-at-a-time replies.

That 24-hour period is short enough to not cause serious disruption of legitimate traffic, since real SMTP implementations will keep trying to deliver for at least a few days. Experience from large-scale implementations of the technique shows that it rarely produces false positives. Machines that continue spamming after 24 hours will make it back to the tarpit soon enough.

## Setting Up a Traplist

To set up your traplist, use `spamdb`'s `-T` option. In my case, the strange address<sup>4</sup> I mentioned earlier in “Greylisting in Practice” on page 96 was a natural candidate for inclusion:

---

```
$ sudo spamdb -T -a wkitp98zpu.fsf@datadok.no
```

---

The command I actually entered was `$ sudo spamdb -T -a "<wkitp98zpu.fsf@datadok.no>"`. In OpenBSD 4.1 and newer, `spamdb` does not require the angle brackets or quotes, but it will accept them.

Add as many addresses as you like. I tend to find new additions for my local list of spamtrap addresses by looking in the greylist and mail server logs for failed attempts to deliver delivery failure reports to nonexistent addresses in my domain (yes, it really is as crazy as it sounds).

4. Of course, this address is totally bogus. It looks like the kind of message ID the GNUS email and news client generates, and it was probably lifted from a news spool or some unfortunate malware victim’s mailbox.

**WARNING** Make sure that the addresses you add to your *spamtrap* lists are invalid and will stay invalid. There is nothing quite like the embarrassment of discovering that you made a valid address into a *spamtrap*, however temporarily.

The following log fragment shows how a spam-sending machine is grey-listed at first contact, and then comes back and clumsily tries to deliver messages to the curious address I added to my traplist, only to end up in the *spamd-greytrap* blacklist after a few minutes. We know what it will be doing for the next 20-odd hours.

---

```
Nov  6 09:50:25 delilah spamd[23576]: 210.214.12.57: connected (1/0)
Nov  6 09:50:32 delilah spamd[23576]: 210.214.12.57: connected (2/0)
Nov  6 09:50:40 delilah spamd[23576]: (GREY) 210.214.12.57: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
Nov  6 09:50:40 delilah spamd[23576]: 210.214.12.57: disconnected after 15 seconds.
Nov  6 09:50:42 delilah spamd[23576]: 210.214.12.57: connected (2/0)
Nov  6 09:50:45 delilah spamd[23576]: (GREY) 210.214.12.57: <bounce-3C7E40A4B3@branch15.summer-
bargainz.com> -> <adm@dataped.no>
Nov  6 09:50:45 delilah spamd[23576]: 210.214.12.57: disconnected after 13 seconds.
Nov  6 09:50:50 delilah spamd[23576]: 210.214.12.57: connected (2/0)
Nov  6 09:51:00 delilah spamd[23576]: (GREY) 210.214.12.57: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
Nov  6 09:51:00 delilah spamd[23576]: 210.214.12.57: disconnected after 18 seconds.
Nov  6 09:51:02 delilah spamd[23576]: 210.214.12.57: connected (2/0)
Nov  6 09:51:02 delilah spamd[23576]: 210.214.12.57: disconnected after 12 seconds.
Nov  6 09:51:02 delilah spamd[23576]: 210.214.12.57: connected (2/0)
Nov  6 09:51:18 delilah spamd[23576]: (GREY) 210.214.12.57: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
Nov  6 09:51:18 delilah spamd[23576]: 210.214.12.57: disconnected after 16 seconds.
Nov  6 09:51:18 delilah spamd[23576]: (GREY) 210.214.12.57: <bounce-3C7E40A4B3@branch15.summer-
bargainz.com> -> <adm@dataped.no>
Nov  6 09:51:18 delilah spamd[23576]: 210.214.12.57: disconnected after 16 seconds.
Nov  6 09:51:20 delilah spamd[23576]: 210.214.12.57: connected (1/1), lists: spamd-greytrap
Nov  6 09:51:23 delilah spamd[23576]: 210.214.12.57: connected (2/2), lists: spamd-greytrap
Nov  6 09:55:33 delilah spamd[23576]: (BLACK) 210.214.12.57: <gilbert@keyholes.net> ->
<wkitp98zpu.fsf@datadok.no>
Nov  6 09:55:34 delilah spamd[23576]: (BLACK) 210.214.12.57: <bounce-3C7E40A4B3@branch15.summer-
bargainz.com> -> <adm@dataped.no>
```

---

As a side note, it looks like even though the spammer moved to send from a different machine, both the *From:* and *To:* addresses stayed the same. The fact that he is still trying to send to an address that has never been deliverable is a strong indicator that this spammer does not check his lists frequently.

### ***Managing Lists with spamdb***

There may be times when you need to view or change the contents of blacklists, whitelists, and greylists. These records are located in the */var/db/spamdb* database, and an administrator's main interface to managing those lists is *spamdb*.

Early versions of `spamdb` simply offered options to add whitelist entries to the database or update existing ones (`spamdb -a nn.mm.nn.mm`). You could delete whitelist entries (`spamdb -d nn.mm.nn.mm`) to compensate for shortcomings in either the blacklists used or the effects of the greylisting algorithms. Recent versions of `spamdb` offer some interesting features to support greytrapping.

### Updating Lists

If you run `spamdb` without any parameters, it lists the contents of your `spamdb` database, and it lets you add or delete both spamtrap addresses and traplist entries. You can also add whitelist entries on the fly.

If you want to add a host to your whitelist without adding it to your permanent *nospamd* file and reloading your rule set or the table, you could do it from the command line instead, like this:

---

```
$ sudo spamdb -a 213.187.179.198
```

---

If a spam sender managed to get a message delivered despite your best efforts, you could correct the situation by adding the spam sender to the `spamd-greytrap` list like this:

---

```
$ sudo spamdb -a -t 192.168.2.128
```

---

Adding a new trap address is just as simple:

---

```
$ sudo spamdb -a -T _medvetsky@ehtrib.org
```

---

If you want to reverse either of these decisions, you would simply substitute `-d` for the `-a` option in both these commands.

### Keeping `spamd` Greylists in Sync

Beginning with OpenBSD 4.1, `spamd` can keep greylisting databases in sync across any number of cooperating greylisting gateways. The implementation is via a set of `spamd` command-line options:

- The `-Y` option specifies a *sync target*—that is, the IP address(es) of other `spamd` running gateways you want to inform of updates to your greylisting information.
- On the receiving end, the `-y` option specifies a *sync listener*, which is the address or interface where this `spamd` instance is prepared to receive greylisting updates from other hosts.

For example, our main `spamd` gateway `mainoffice-gw.example.com` might have the following options added to its startup command line to establish a sync target and sync listener, respectively:

---

```
-Y minorbranch-gw.example.com -y mainoffice-gw.example.com
```

---

Conversely, `minorbranch-gw.example.com` at the branch office would have the hostnames reversed:

---

```
-Y mainoffice-gw.example.com -y minorbranch-gw.example.com
```

---

The `spamd` daemon also supports shared-secret authentication between the synchronization partners. Specifically, if you create the file `/etc/mail/spamd.key` and distribute copies of it to all synchronization partners, it will be used to calculate the necessary checksums for authentication. The `spamd.key` file itself can be any kind of data, such as random data harvested from `/dev/arandom` as suggested by the `spamd` man page.

**NOTE** *In situations where direct synchronization of `spamd`-related data is not practical, or if you simply want to share your `spamd`-`greytrap` with others, exporting the contents of your list of locally trapped spam senders to a text file may be desirable. The list format `spamd-setup` expects is one address per line, optionally with comment lines starting with one or more # characters. Exporting your list of currently trapped addresses in a usable format can be as simple as putting together a one-liner with `spamdb`, `grep`, and a little imagination.*

### **Detecting Out-of-Order MX Use**

OpenBSD 4.1 gave `spamd` the ability to detect out-of-order MX use. Contacting a secondary mail exchanger first instead of trying the main one is a fairly well-known spammer trick, and one that runs contrary to the behavior we expect from ordinary email transfer agents. In other words, if someone tries the email exchangers in the wrong order, we can be pretty sure that they're trying to deliver spam.

For our `example.com` domain with main mail server 192.0.2.225 and backup 192.0.2.224, adding `-M 192.0.2.224` to `spamd`'s startup options would mean that any host that tries to contact 192.0.2.224 via SMTP before contacting the main mail server at 192.0.2.225 will be added to the local `spamd`-`greytrap` list for the next 24 hours.

### **Handling Sites That Do Not Play Well with Greylisting**

Unfortunately, there are situations where you will need to compensate for the peculiarities of other sites' email setups.

The first email message sent from any site that has not contacted you for as long as the greylist keeps its data around will be delayed for some random amount of time, which depends mainly on the sender's retry interval. There are times when even a minimal delay is undesirable. If, for example, you have some infrequent customers who always demand your immediate and urgent attention to their business when they do contact you, an initial delivery delay of what could be up to several hours may not be optimal. In addition, you are bound to encounter misconfigured mail servers that either do not retry at all or retry too quickly, perhaps stopping delivery retries after just one attempt.

Also, some sites are large enough to have several outgoing SMTP servers, and they do not play well with greylisting since they are not guaranteed to retry delivery of any given message from the same IP address used with the prior delivery attempt. Even though those sites comply with the retry requirements, it's obvious that this is one of the few remaining downsides of greylisting.

One way to compensate for such situations is to define a table for a local whitelist to be fed from a file in case of reboots. To make sure SMTP traffic from the addresses in the table is not fed to `spamd`, add a `pass` rule to allow the traffic to pass:

---

```
table <nospamd> persist file "/etc/mail/nospamd"
pass in log on egress proto tcp from <nospamd> to port smtp
```

---

In pre-OpenBSD 4.7 syntax, add a `no rdr` rule at the top of your redirection block and a matching `pass` rule to let SMTP traffic from the hosts in your `nospamd` table through, as shown here:

---

```
no rdr proto tcp from <nospamd> to $mailservers port smtp
pass in log on egress proto tcp from <nospamd> to port smtp
```

---

Once you have made these changes to your rule set, enter the addresses you need to protect from redirection into the `/etc/mail/nospamd` file, and then reload your rule set using `pfctl -f /etc/pf.conf`. You can then use all the expected table tricks on the `<nospamd>` table, including replacing its content after editing the `nospamd` file. In fact, this approach is strongly hinted at in both man pages and sample configuration files distributed with recent versions of `spamd`.

At least some sites with many outgoing SMTP servers publish information about which hosts are allowed to send email for their domain via Sender Policy Framework (SPF) records as part of the domain's DNS information.<sup>5</sup> To retrieve the SPF records for our `example.com` domain, use the `host` command's `-ttxt` option like this:

---

```
$ host -ttxt example.com
```

---

The command would produce an answer roughly like this:

---

```
example.com descriptive text "v=spf1 ip4:192.0.2.128/25 -all"
```

---

Here, the text in quotes is the `example.com` domain's SPF record. If you want email from `example.com` to arrive quickly, and you trust the people there not to send or relay spam, choose the address range from the SPF record, add it to your `nospamd` file, and reload the `<nospamd>` table contents from the updated file.

---

5. SPF records are stored in DNS zones as TXT records. See <http://www.openspf.org/> for details.

## Spam-Fighting Tips

When used selectively, blacklists combined with `spamd` are powerful, precise, and efficient spam-fighting tools. The load on the `spamd` machine is minimal. On the other hand, `spamd` will never perform better than its weakest data source, which means you will need to monitor your logs and use whitelisting when necessary.

It is also feasible to run `spamd` in a pure greylisting mode, with no blacklists. In fact, some users report that a purely greylisting `spamd` configuration is about as effective a spam-fighting tool as configurations with blacklists, and sometimes significantly more effective than content filtering. One such report that was posted to *openbsd-misc* (accessible among other places via <http://marc.info/>; search for the subject “Followup - `spamd` greylisting results”) claimed that a pure greylisting configuration immediately rid the company of approximately 95 percent of its spam load.

I recommend two very good blacklists. One is Bob Beck’s “ghosts of usenet postings past”–based traplist, generated automatically by computers running `spamd` at the University of Alberta. Bob’s setup is a regular `spamd` system that removes trapped addresses automatically after 24 hours, which means that you get an extremely low number of false positives. The number of hosts varies widely and has been as high as 670,000. While still officially in testing, the list was made public in January 2006. The list is available from <http://www.openbsd.org/spamd/traplist.gz>. It is part of recent sample `spamd.conf` files as the `uatraps` blacklist.

The other list I recommend is *heise.de*’s `nixspam`, which has a 12-hour automatic expiry and extremely good accuracy. It’s also in the sample `spamd.conf` file. Detailed information about this list is available from [http://www.heise.de/ix/nixspam/dnsbl\\_en/](http://www.heise.de/ix/nixspam/dnsbl_en/).

Once you’re happy with your setup, try introducing local greytrapping. This is likely to catch a few more undesirables, and it’s good, clean fun. Some limited experiments, carried out while writing this chapter (chronicled at <http://bsdly.blogspot.com/>, entries starting with <http://bsdly.blogspot.com/2007/07/hey-spammer-heres-list-for-you.html>), even suggest that harvesting the invalid addresses spammers use from your mail server logs, from `spamd` logs, or directly from your greylist to put in your traplist is extremely efficient. Publishing the list on a moderately visible web page appears to ensure that the addresses you put there will be recorded over and over again by address-harvesting robots, and will provide you with even better greytrapping material, since they are then more likely to be kept on the spammers’ list of known good addresses.

# 7

## QUEUES, SHAPING, AND REDUNDANCY



This chapter deals with two main topics, which taken separately or together have the potential to radically transform your networking experience. The common theme in this chapter is managing resource availability.

In the first part of the chapter, we look at how to use the ALTQ traffic-shaping subsystem to allocate bandwidth resources efficiently and according to a specified policy. The second part of the chapter explores how to make sure your resources stay available by using the redundancy features offered by the CARP and pfsync protocols.

### Directing Traffic with ALTQ

*ALTQ*, short for *ALternate Queueing*, is a very flexible mechanism for network traffic shaping, which lived a life of its own before it was integrated into PF on OpenBSD.<sup>1</sup> On OpenBSD, ALTQ was integrated into the PF code for the

<sup>1</sup>. The original research on ALTQ was presented in a paper for the USENIX 1999 conference. You can read Kenjiro Cho's paper "Managing Traffic with ALTQ" online at <http://www.usenix.org/publications/library/proceedings/usenix99/cho.html>.

OpenBSD 3.3 release by Henning Brauer, with the configuration done in *pf.conf* mainly because it makes sense to integrate traffic shaping and filtering. PF ports to other BSDs followed suit, with at least some optional ALTQ integration.

## **Basic ALTQ Concepts**

Managing your bandwidth has a lot in common with balancing your checkbook or handling other resources that are either scarce or available in finite quantities. The resource is available in a constant supply with hard upper limits, and you need to allocate the resource with maximum *efficiency*, according to the *priorities* set out in your *policy* or *specification*.

The core of ALTQ bandwidth management is the *queue* concept. Queues are a form of buffers for network packets. The packets are held in a queue until they are either dropped or sent on their way according to the criteria that apply to the queue, and subject to the queue's available bandwidth. Queues are attached to specific interfaces, and bandwidth is managed on a per-interface basis, with available bandwidth on a given interface subdivided into the queues you define.

Queues are defined with a specific amount of bandwidth, a specific portion of available bandwidth, or sometimes hierarchical priority. *Priority* in this context is an indicator of preference or which queue should be serviced with the shortest delay. Some queue types can even be configured with a combination of bandwidth allocation and priority. For even further refinement, with some queue types, you can allocate portions of each queue's bandwidth share to subqueues, or queues within queues, which share the parent queue's resources.

Once queue definitions are in place, you integrate traffic shaping into your rule set by rewriting your `pass` or `match` rules to assign traffic to a specific queue.

**NOTE** *In ALTQ setups, any traffic that you do not explicitly assign to a specific queue gets lumped in with everything else in the default queue.*

## **Queue Schedulers, aka Queue Disciplines**

In the default networking setup, with no ALTQ-style queueing, the TCP/IP stack and its filtering subsystem process the packets in order as they arrive on an interface. This is what we generally refer to as the *first in, first out* (FIFO) discipline.

ALTQ queues can be set up to behave quite differently, sometimes with startling effect. Each of the three queue scheduler algorithms, or *disciplines*, offer their own unique sets of options. The types are `priq`, `cbq`, and `hfsc`.

### **priq**

*Priority-based queues* are defined purely in terms of priority within the total bandwidth. For `priq` queues, the allowed priority range is 0 through 15, where a higher value earns preferential treatment. Packets that match the

criteria for higher priority queues are serviced before the ones matching lower priority queues.

### **cbq**

*Class-based queues* are defined as constant-size bandwidth allocations, as a percentage of the total available or in units of kilobits, megabits, or gigabits per second. A cbq queue can be subdivided into queues that are also assigned priorities in the range 0 to 7, and again a higher priority means preferential treatment.

Packets are kept in the queue until the bandwidth is available. For queues that are subdivided into queues with priority as well as bandwidth allocations, packets that match the criteria for a higher priority queue are serviced sooner.

### **hfsc**

The hfsc discipline uses the *Hierarchical Fair Service Curve* (HFSC) algorithm to ensure a “fair” allocation of bandwidth among the queues in a hierarchy. HFSC comes with the possibility of setting up queueing regimes with guaranteed minimum allocations and hard upper limits. You can even have allocations that vary over time, as well as fine-grained priority with a 0 through 7 range.

Both the algorithm and the corresponding setup are fairly complicated, with a number of tunable parameters. For that reason, most ALTQ practitioners stick with the simpler queue types. Yet the ones who claim to understand HFSC pretty much swear by it.

## **Setting Up ALTQ**

Enabling ALTQ so you can use the queueing logic in your rule sets may require some extra steps, depending on which operating system is your platform of choice.

### **ALTQ on OpenBSD**

On OpenBSD, all supported queue disciplines are compiled into the GENERIC and GENERIC.MP kernels. The only configuration you need to do involves editing your *pf.conf*.

### **ALTQ on FreeBSD**

On FreeBSD, you need to check that your kernel has ALTQ and the ALTQ queue discipline options compiled in. The default FreeBSD GENERIC kernel does not have ALTQ options enabled, as you may have noticed from the messages you saw when running the */etc/rc.d/pf* script to enable PF. The relevant options are listed here:

---

options	ALTQ	
options	ALTQ_CBQ	# Class Bases Queueing (CBQ)
options	ALTQ_RED	# Random Early Detection (RED)
options	ALTQ_RIO	# RED In/Out
options	ALTQ_HFSC	# Hierarchical Packet Scheduler (HFSC)

---

```
options      ALTO_PRIQ      # Priority Queueing (PRIQ)
options      ALTO_NOPCC     # Required for SMP build
```

---

The `ALTO` option is needed to enable ALTQ in the kernel. On SMP systems, you also need the `ALTO_NOPCC` option. Depending on which types of queues you will be using, you need to enable at least of one of `ALTO_CBQ`, `ALTO_PRIQ`, or `ALTO_HFSC`. Finally, you can enable the congestion-avoidance techniques *Random Early Detection* (RED) and *RED In/Out* with the `ALTO_RED` and `ALTO_RIO` options, respectively. See the *FreeBSD Handbook* for information about how to compile and install a custom kernel with these options.

### ALTQ on NetBSD

ALTQ was integrated in the NetBSD 4.0 PF implementation and is supported in NetBSD 4.0 and later releases. NetBSD's default GENERIC kernel configuration does not include the ALTQ-related options. However, the GENERIC configuration file comes with all relevant options commented out for easy inclusion. The following are the main kernel options:

---

```
options      ALTO          # Manipulate network interfaces' output queues
options      ALTO_CBQ       # Class-Based Queueing
options      ALTO_HFSC     # Hierarchical Fair Service Curve
options      ALTO_PRIQ      # Priority Queueing
options      ALTO_RED       # Random Early Detection
```

---

The `ALTO` option is needed to enable ALTQ in the kernel. Depending on which types of queues you will be using, you need to enable at least of one of `ALTO_CBQ`, `ALTO_PRIQ`, or `ALTO_HFSC`.

Using ALTQ requires that you compile PF into the kernel. The PF loadable module does not support the ALTQ functionality. Check the NetBSD PF documentation at <http://www.netbsd.org/Documentation/network/pf.html> for the most up-to-date information.

## Setting Up Queues

Now we will look at examples of setting up the three different types of ALTQ queues: priority-based, class-based bandwidth, and HFSC.

The general syntax for ALTQ queues in PF looks like this:

---

```
altq on interface type [options ...] main_queue { sub_q1, sub_q2 ...}
    queue sub_q1 [ options ... ]
    queue sub_q2 [ options ... ] { subA, subB, ... }
[...]
pass [ ... ] queue sub_q1
pass [ ... ] queue sub_q2
```

---

Note that cbq and hfsc queues can have several levels of subqueues. The priq queues are essentially flat, with only one queue level. We will address syntax specifics for each type in the following sections.

## WHAT IS YOUR TOTAL USABLE BANDWIDTH?

It can be difficult to determine actual usable bandwidth on a specific interface for queuing. If you do not specify a total bandwidth, the total bandwidth available will be used to calculate the allocations. However, some types of interfaces cannot reliably report the actual bandwidth value. One common example of this discrepancy is where your gateway's external interface is a 100 megabit (Mb) Ethernet interface, attached to a DSL line that actually offers only 8Mb download and 1Mb upload.\* The Ethernet interface will then confidently report 100Mb bandwidth, not the DSL values.

Therefore, it usually makes sense to set the total bandwidth to a fixed value. Unfortunately, the value to use may not be exactly what your bandwidth supplier tells you is available. There will always be some overhead that varies slightly over the various technologies and implementations. In typical TCP/IP over wired Ethernet, overhead can be as low as single-digit percentages. TCP/IP over ATM has been known to have overhead as high as almost 20 percent. If your bandwidth supplier does not provide the overhead information, you will need to make an educated guess at the start value for your experimentation. In any case, you should be acutely aware that the total bandwidth available is never larger than the bandwidth of the weakest link in your network path.

Queues are supported only for outbound connections relative to the system doing the queueing. When planning your bandwidth management, you should consider the actual usable bandwidth to be equal to the weakest (lowest bandwidth) link in the connection's path, even if your queues are set up on a different interface.

\* This really dates the book, I know. In a few years, these numbers will seem quaint.

## Priority-Based Queues

The basic concept for priority-based queues (`priq`) is fairly straightforward and perhaps the easiest to understand. Within the total bandwidth allocated to the main queue, all that matters is traffic priority. You assign queues a priority value in the range 0 through 15, where a higher value means that the queue's requests for traffic are serviced sooner.

### A Real-World Example

For a real-world example, we can look to Daniel Hartmeier. He discovered a simple yet effective way to improve the throughput for his home network by using ALTQ. Like many people, Daniel had his home network on an asymmetric connection, with total usable bandwidth low enough that he felt a strong desire to get better bandwidth utilization. In addition, when the line was running at or near capacity, some oddities started appearing. One symptom in particular seemed to indicate there was room for improvement: Incoming traffic (downloads, incoming mail, and such) slowed down disproportionately whenever outgoing traffic started—more than could be explained by measuring the raw amount of data transferred. It all came back to a basic feature of TCP.

When a TCP packet is sent, the sender expects acknowledgment (in the form of an ACK packet) from the receiving end and will wait for a specified time for it to arrive. If the ACK does not arrive within the specified time, the

sender assumes that the packet has not been received and resends the packet it originally sent. The problem is that in a default setup, packets are serviced sequentially by the interface as they arrive. This inevitably means that ACK packets, with essentially no data payload, wait in line while the larger data packets are transferred.

A testable hypothesis formed: If the tiny, practically data-free ACK packets were able to slip in between the larger data packets, this would lead to a more efficient use of available bandwidth. The simplest practical way to implement and test the theory was to set up two queues with different priorities and integrate them into the rule set. The following lines show the relevant parts of the rule set.

---

```
ext_if="kueo"

altq on $ext_if priq bandwidth 100Kb queue { q_pri, q_def }
    queue q_pri priority 7
    queue q_def priority 1 priq(default)

pass out on $ext_if proto tcp from $ext_if queue (q_def, q_pri)

pass in on $ext_if proto tcp to $ext_if queue (q_def, q_pri)
```

---

Here, the priority-based queue is set up on the external interface, with two subordinate queues. The first subqueue, `q_pri`, has a high priority value of 7; the other subqueue (`q_def`) has a significantly lower priority value of 1.

This seemingly simple rule set works by exploiting how ALTQ treats queues with different priorities. Once a connection is assigned to the main queue, ALTQ inspects each packet's type of service (ToS) field. ACK packets have the ToS delay bit set to low, which indicates that the sender wanted the speediest delivery possible. When ALTQ sees a low-delay packet, and queues of differing priorities are available, it will assign the packet to the higher priority queue. This means that the ACK packets skip ahead of the lower priority queue and are delivered more quickly, which in turn means that data packets are serviced more quickly.

The net result is that a configuration like this provides better performance than a pure FIFO configuration with the same hardware and available bandwidth.<sup>2</sup>

### Using a match Rule for Queue Assignment

In the preceding example, the rule set was constructed the traditional way, with the queue assignment as part of the `pass` rules. However, this is not the only way to do queue assignment. Using `match` rules (available in OpenBSD 4.6 and later), it is incredibly easy to retrofit this simple priority-queueing regime onto an existing rule set.

---

2. Daniel's article about this version of his setup, at <http://www.benzedrine.cx/ackpri.html>, contains a more detailed analysis.

If you have been working through the examples in Chapters 3 and 4, it is likely that your rule set already has a `match` rule that applies `nat-to` on your outgoing traffic. If you want to introduce priority-based queueing to your rule set, you can get started by adding the queue definitions and making some minor adjustments to your outgoing `match` rule.

Start with the queue definition from the preceding example, and adjust the total bandwidth to local conditions, as in this example:

---

```
altq on $ext_if priq bandwidth $ext_bw queue { q_pri, q_def }
    queue q_pri priority 7
    queue q_def priority 1 priq(default)
```

---

This gives the queues whatever bandwidth allocation you define with the `ext_bw` macro.

The simplest and quickest way to integrate the queues into your rule set is to edit your outgoing `match` rule to read something like this:

---

```
match out on $ext_if from $int_if:network nat-to ($ext_if) queue (q_def, q_pri)
```

---

Reload your rule set, and the priority-queueing regime is applied to all traffic that is initiated from your local network.

You can use the `systat` command to get a live view of how traffic is assigned to your queues.

---

```
$ sudo systat queues
```

---

This will give you a live display something like this:

---

2 users	Load	0.39	0.27	0.30		Fri	Apr	1	16:33:44	2011		
QUEUE	BW	SCH	PR	PKTS	BYTES	DROP_P	DROP_B	QLEN	BORRO	SUSPE	P/S	B/S
q_pri		priq	7	21705	1392K	0	0	0		12	803	
q_def		priq		12138	6759K	0	0	0		9	4620	

---

Looking at the numbers in the `PKTS` (packets) and `BYTES` columns, you see a clear indication that the queuing is working as intended.

The `q_pri` queue has processed a rather large number of packets in relation to the amount of data, just as we expected. The ACK packets do not take up a lot of space. On the other hand, the traffic assigned to the `q_def` queue has more data in each packet, and the numbers show essentially the reverse packet numbers-to-data size ratio compared to the `q_pri` queue.

**NOTE** *systat* is a rather capable program on all BSDs, and the OpenBSD version offers several views that are relevant to PF, so far not found in the *systat* variants on the other systems. We will be looking at *systat* again in the next chapter. In the meantime, read the man pages and play with the program. It's a very useful tool for getting to know your system.

## **Class-Based Bandwidth Allocation for Small Networks**

Maximizing network performance generally feels nice. However, you may find that your network has other needs. For example, it might be important for some traffic—such as mail and other vital services—to have a baseline amount of bandwidth available at all times, while other services—peer-to-peer file sharing comes to mind—should not be allowed to consume more than a certain amount. For addressing these kinds of requirements or concerns, the class-based queue (cbq) discipline offers a slightly larger set of options.

To illustrate how to use cbq, we move on to another example, which builds on the rule sets from previous chapters. This is a small local network, where we want to let the users on the local network connect to a predefined set of services outside their own network and also access a web server somewhere on the local network.

### **Queue Definition**

Here, all queues are set up on the external, Internet-facing interface. This approach makes sense mainly because bandwidth is more likely to be limited on the external link than on the local network. In principle, however, allocating queues and running traffic shaping can be done on any network interface. The following example setup includes a cbq queue for a total bandwidth of 2Mb with six subqueues.

---

```
altq on $ext_if cbq bandwidth 2Mb queue { main, ftp, udp, web, ssh, icmp }
    queue main bandwidth 18% cbq(default borrow red)
    queue ftp bandwidth 10% cbq(borrow red)
    queue udp bandwidth 30% cbq(borrow red)
    queue web bandwidth 20% cbq(borrow red)
    queue ssh bandwidth 20% cbq(borrow red) { ssh_interactive, ssh_bulk }
        queue ssh_interactive priority 7 bandwidth 20%
        queue ssh_bulk priority 0 bandwidth 80%
    queue icmp bandwidth 2% cbq
```

---

The subqueue `main` has 18 percent of the bandwidth and is designated as the default queue. This means any traffic that matches a pass rule but is not explicitly assigned to some other queue ends up here. The `borrow` and `red` keywords mean that the queue may “borrow” bandwidth from its parent queue, while the system attempts to avoid congestion by applying the RED algorithm.

The other queues follow more or less the same pattern, up to the subqueue `ssh`, which itself has two subqueues with separate priorities. Here we see a variation on the ACK priority example. Bulk SSH transfers, typically SCP file transfers, are transmitted with a ToS indicating throughput, while interactive SSH traffic has the ToS flag set to low delay and skips ahead of the bulk transfers. The interactive traffic is likely to be less bandwidth-consuming and gets a smaller share of the bandwidth, but receives preferential treatment because of the higher priority value assigned to it.

This scheme also helps the speed of SCP file transfers, since the ACK packets for the SCP transfers will be assigned to the higher priority subqueue.

Finally, we have the `icmp` queue, which is reserved for the remaining 2 percent of the bandwidth, from the top level. This guarantees a minimum amount of bandwidth for ICMP traffic that we want to pass but that does not match the criteria for being assigned to the other queues.

### Rule Set

To make it all happen, we use these `pass` rules, which indicate which traffic is assigned to the queues and their criteria:

---

```
set skip on { lo, $int_if }
pass log quick on $ext_if proto tcp to port ssh queue (ssh_bulk, ssh_interactive)
pass in quick on $ext_if proto tcp to port ftp queue ftp
pass in quick on $ext_if proto tcp to port www queue http
pass out on $ext_if proto udp queue udp
pass out on $ext_if proto icmp queue icmp
pass out on $ext_if proto tcp from $localnet to port $client_out
```

---

The rules for `ssh`, `ftp`, `www`, `udp`, and `icmp` assign traffic to their respective queues. The last catchall rule passes all other traffic from the local network, lumping it into the default `main` queue.

### A Basic HFSC Traffic Shaper

The simple schedulers we have looked at so far can make for rather efficient setups, but network admins with traffic-shaping ambitions tend to look for a little more flexibility than can be found in the pure priority-based queues or the simple class-based variety. After the gentle queuing introduction we've just been through, how does it sound to you if I say there is a regime with flexible bandwidth allocation, guaranteed lower and upper bounds for bandwidth available to each queue, and variable allocations over time—and one that only starts shaping when there is an actual need?

The HFSC queueing algorithm (`hfsc` in `pf.conf` terminology) offers all of these, but the added flexibility comes at a price: The setup is a tad more complex than the other types, and tuning your setup for an optimal result can be quite an interesting process.

### Queue Definition

Working from the same configuration we altered slightly earlier, we insert this queue definition early in the `pf.conf` file:

---

```
altq on $ext_if bandwidth $ext_bw hfsc queue { main, spamd }
queue main bandwidth 99% priority 7 qlimit 100 hfsc (realtime 20%, linkshare 99%) \
    { q_pri, q_def, q_web, q_dns }
queue q_pri bandwidth 3% priority 7 hfsc (realtime 0, linkshare 3% red )
queue q_def bandwidth 47% priority 1 hfsc (default realtime 30% linkshare 47% red)
queue q_web bandwidth 47% priority 1 hfsc (realtime 30% linkshare 47% red)
queue q_dns bandwidth 3% priority 7 qlimit 100 hfsc (realtime (30Kb 3000 12Kb), \
    linkshare 3%)
queue spamd bandwidth 0% priority 0 qlimit 300 hfsc (realtime 0, upperlimit 1%, \
    linkshare 1%)
```

---

As you can see, the `hfsc` queue definitions take slightly different parameters than the simpler disciplines. We start off with this rather small hierarchy by splitting the top-level queue into two. At the next level, we subdivide the `main` queue into several subqueues, each with a defined priority. All the subqueues have a `realtime` value set; this is the guaranteed minimum bandwidth allocated to the queue. The optional `upperlimit` sets a hard upper limit on the queue's allocation. The `linkshare` parameter sets the allocation the queue will have available when it is backlogged; that is, it has started to eat into its `qlimit` allocation.

In case of congestion, each queue by default has a pool of 50 slots, the queue limit (`qlimit`), to keep packets around when they cannot be transmitted immediately. In this example, the top-level queues `main` and `spamd` both have larger than default pools set by their `qlimit` setting: 100 for `main` and 300 for `spamd`. Cranking up queue sizes here means we are a little less likely to drop packets when the traffic approaches the set limits, but it also means that when the traffic shaping kicks in, we will see increased latency for connections that end up in these larger than default pools.

The queue hierarchy here uses two familiar tricks to make efficient use of available bandwidth:

- It uses a variation of the high- and low-priority mix demonstrated in the earlier pure priority example.
- We speed up almost all other traffic (and most certainly the web traffic that appears to be the main priority here) by allocating a small but guaranteed portion of bandwidth for name service lookups. For the `q_dns` queue, we set up the `realtime` value with a time limit—after 3000 milliseconds, the `realtime` allocation goes down to 12kb. This can be useful to speed connections that transfer most of their payload during the early phases.

## Rule Set

The next step is to tie the newly created queues into the rule set. Assuming you have a filtering regime in place already, the tie-in becomes amazingly simple by adding a few `match` rules:

---

```
match out on $ext_if from $air_if:network nat-to ($ext_if) queue (q_def, q_pri)
match out on $ext_if from $int_if:network nat-to ($ext_if) queue (q_def, q_pri)
match out on $ext_if proto tcp to port { www https } queue (q_web, q_pri)
match out on $ext_if proto { tcp udp } to port domain queue (q_dns, q_pri)
match out on $ext_if proto icmp queue (q_dns, q_pri)
```

---

Here, the `match` rules once again do the ACK packet speedup trick with the high- and low-priority queue assignment, just as you saw earlier in the pure priority-based system. The only exception is when we assign traffic to our lowest priority queue, where we really do not care to have any speedup at all:

---

```
pass in log on egress proto tcp to port smtp rdr-to 127.0.0.1 port spamd queue spamd
```

---

This has the intention of slowing down the spammers a little more on their way to our `spamd`. With a hierarchical queue system in place, `systat` queues shows the queues and their traffic as a hierarchy, too.

2 users	Load	0.22	0.25	0.25	Fri Apr 1 16:43:37 2011							
QUEUE	BW	SCH	PRI0	PKTS	BYTES	DROP_P	DROP_B	QLEN	BORROW	SUSPEN	P/S	B/S
root_nfe0	20M	hfsc	0	0	0	0	0	0	0	0	0	0
main	19M	hfsc	7	0	0	0	0	0	0	0	0	0
q_pri	594K	hfsc	7	1360	82284	0	0	0	0	11	770	
q_def	9306K	hfsc		158	15816	0	0	0	0	0.2	11	
q_web	9306K	hfsc		914	709845	0	0	0	0	50	61010	
q_dns	594K	hfsc	7	196	17494	0	0	0	0	3	277	
spamd	0	hfsc	0	431	24159	0	0	0	0	2	174	

The root queue is shown as attached to the physical interface—as `nfe0` and `root_nfe0` in this case. `main` and its subqueues `q_pri`, `q_def`, `q_web`, and `q_dns` are shown with their bandwidth allocations and number of bytes and packets passed. The `DROP_P` and `DROP_B` columns are where number of packets and bytes dropped, respectively, would appear if we had been forced to drop packets at this stage. The final two columns show live updates of packets per second and bytes per second.

### ***Queueing for Servers in a DMZ***

Back in Chapter 5, we set up a network with a single gateway but all externally visible services configured on a separate DMZ network. That way, all traffic to the servers from both the Internet and the internal network has to pass through the gateway.

The network schematic is shown in Figure 7-1, which is identical to Figure 5-2.

With the rule set from Chapter 5 as the starting point, we will be adding some queuing in order to optimize our network resources. The physical and logical layout of the network will not change.

The most likely bottleneck for this network is the bandwidth for the connection between the gateway’s external interface and the Internet at large. The bandwidth elsewhere in our setup is not infinite, of course, but the available bandwidth on any interface in the local network is likely to be less of a limiting factor than the bandwidth actually available for communication with the outside world. For services to be available with the best possible performance, we need to set up the queues so the bandwidth available at the site is made available to the traffic we want to allow.

In our example, it is likely that the interface bandwidth on the DMZ interface is either 100Mb or 1Gb, while the *actual available bandwidth* for connections from outside the local network is considerably smaller. This consideration shows up in our queue definitions, where you clearly see that the actual bandwidth available for external traffic is the main limitation in the queue setup.

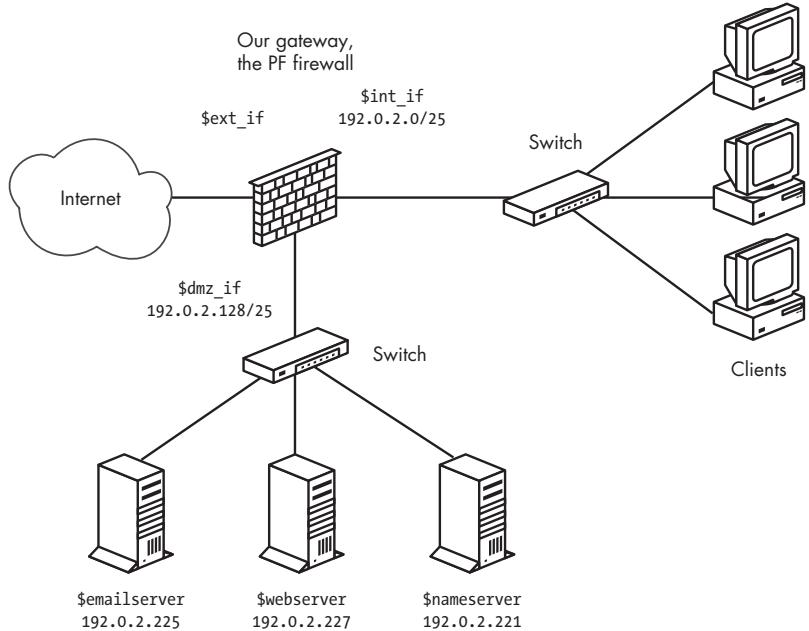


Figure 7-1: Network with DMZ

---

```

total_ext = 2Mb
total_dmz = 100Mb
altq on $ext_if cbq bandwidth $total_ext queue { ext_main, ext_web, ext_udp, \
    ext_mail, ext_ssh }
queue ext_main bandwidth 25% cbq(default borrow red) { ext_hi, ext_lo }
queue ext_hi priority 7 bandwidth 20%
queue ext_lo priority 0 bandwidth 80%
queue ext_web bandwidth 25% cbq(borrow red)
queue ext_udp bandwidth 20% cbq(borrow red)
queue ext_mail bandwidth 30% cbq(borrow red)
altq on $dmz_if cbq bandwidth $total_dmz queue { ext_dmz, dmz_main, dmz_web, \
    dmz_udp, dmz_mail }
queue ext_dmz bandwidth $total_ext cbq(borrow red) queue { ext_dmz_web, \
    ext_dmz_udp, ext_dmz_mail }
    queue ext_dmz_web bandwidth 40% priority 5
    queue ext_dmz_udp bandwidth 10% priority 7
    queue ext_dmz_mail bandwidth 50% priority 3
queue dmz_main bandwidth 25Mb cbq(default borrow red) queue { dmz_main_hi, \
    dmz_main_lo }
queue dmz_main_hi priority 7 bandwidth 20%
queue dmz_main_lo priority 0 bandwidth 80%
queue dmz_web bandwidth 25Mb cbq(borrow red)
queue dmz_udp bandwidth 20Mb cbq(borrow red)
queue dmz_mail bandwidth 20Mb cbq(borrow red)

```

---

Notice that the `total_ext` bandwidth limitation determines the allocation for all queues where the bandwidth for external connections is available. To make use of the new queuing infrastructure, we need to make some changes to the filtering rules, too. Keep in mind that any traffic you do not explicitly assign to a specific queue is assigned to the default queue for the interface. Thus, it is important to tune your filtering rules as well as your queue definitions to the actual traffic in your network.

The main part of the filtering rules could end up looking like this after adding the queues:

---

```
pass in on $ext_if proto { tcp, udp } to $nameservers port domain \
    queue ext_udp
pass in on $int_if proto { tcp, udp } from $localnet to $nameservers \
    port domain
pass out on $dmz_if proto { tcp, udp } to $nameservers port domain \
    queue ext_dmz_udp
pass out on $dmz_if proto { tcp, udp } from $localnet to $nameservers \
    port domain queue dmz_udp
pass in on $ext_if proto tcp to $webserver port $webports queue ext_web
pass in on $int_if proto tcp from $localnet to $webserver port $webports
pass out on $dmz_if proto tcp to $webserver port $webports queue ext_dmz_web
pass out on $dmz_if proto tcp from $localnet to $webserver port $webports \
    queue dmz_web
pass in log on $ext_if proto tcp to $mailserver port smtp
pass in log on $ext_if proto tcp from $localnet to $mailserver port smtp
pass in log on $int_if proto tcp from $localnet to $mailserver port $email
pass out log on $dmz_if proto tcp to $mailserver port smtp queue ext_mail
pass in on $dmz_if from $mailserver to port smtp queue dmz_mail
pass out log on $ext_if proto tcp from $mailserver to port smtp \
    queue ext_dmz_mail
```

---

Notice that only traffic that will pass either the DMZ interface or the external interface is assigned to queues. In this configuration, with no externally accessible services on the internal network, queuing on the internal interface would not make much sense, since it is likely the part of our network with the least restrictions on available bandwidth.

## ***Using ALTQ to Handle Unwanted Traffic***

So far, we have focused on queuing as a method to make sure specific kinds of traffic are let through as efficiently as possible given the conditions that exist in and around your network. Now, we will look at two examples that present a slightly different approach to identify and handle unwanted traffic. These examples demonstrate some queuing-related tricks you can use to keep miscreants in line.

### ***Overloading to a Tiny Queue***

Think back to the “Turning Away the Brutes” section (page 86), where we used a combination of state-tracking options and `overload` rules to fill up a table of addresses for special treatment. The special treatment we

demonstrated in Chapter 6 was to cut all connections, but it is equally possible to assign overload traffic to a specific queue instead.

Consider this rule from our earlier class-based bandwidth example earlier in the chapter:

---

```
pass log quick on $ext_if proto tcp to port ssh flags S/SA \
    keep state queue (ssh_bulk, ssh_interactive)
```

---

We could add state-tracking options, like this:

---

```
pass log quick on $ext_if proto tcp to port ssh flags S/SA \
    keep state (max-src-conn 15, max-src-conn-rate 5/3, \
    overload <bruteforce> flush global) queue (ssh_bulk, ssh_interactive)
```

---

Then we could make one of the queues slightly smaller:

---

```
queue smallpipe bandwidth 1kb cbq
```

---

And then assign traffic from miscreants to the small-bandwidth queue with the following rule:

---

```
pass inet proto tcp from <bruteforce> to port $tcp_services queue smallpipe
```

---

It might also be useful to supplement rules like these with table-entry expiry as described, as described in “Tidying Your Tables with pfctl” on page 89.

### Queue Assignments Based on Operating System Fingerprint

Chapter 6 covered several ways to use `spamd` to cut down on spam. If running `spamd` is not an option in your environment, you can use a queue and rule set based on the common knowledge that machines that send spam are likely to run a particular operating system.

PF has a fairly reliable operating system fingerprinting mechanism, which detects the operating system at the other end of a network connection based on characteristics of the initial SYN packets at connection setup. The following may be a simple substitute for `spamd` if you have determined that legitimate mail is highly unlikely to be delivered from systems that run that particular operating system:

---

```
pass quick proto tcp from any os "Windows" to $ext_if port smtp queue smallpipe
```

---

Here, email traffic originating from hosts that run a particular operating system get no more than 1 kilobit of your bandwidth, with no borrowing.

## **Redundancy and Failover: CARP and pfsync**

High availability and uninterrupted service have been both marketing buzzwords and coveted goals for IT professionals and network administrators as long as most of us can remember. To meet this perceived need and solve a few related problems, CARP and pfsync were added as two highly anticipated features in OpenBSD 3.5.

The Common Address Redundancy Protocol (CARP) was developed as a non-patent-encumbered alternative to the Virtual Router Redundancy Protocol (VRRP), which was quite far along the track to becoming an IETF-sanctioned standard, even though possible patent issues have not been resolved.<sup>3</sup>

One of the main purposes of CARP is to ensure that the network will keep functioning as usual, even when a firewall or other service goes down due to errors or planned maintenance activities such as upgrades.

Complementing CARP, the pfsync protocol is designed to handle synchronization of PF states between redundant packet-filtering nodes or gateways. Both protocols are intended to ensure redundancy for essential network features with automatic failover.

CARP is based on setting up a group of machines as one *master* and one or more redundant *backups*, which are all equipped to handle a common IP address. If the master goes down, one of the backups will inherit the IP address. The handover from one CARP host to another may be authenticated, essentially by setting a shared secret (in practice, much like a password).

In the case of PF firewalls, pfsync can be set up to handle the synchronization, and if the synchronization via pfsync has been properly set up, active connections will be handed over without noticeable interruption. In essence, pfsync is a type of virtual network interface specially designed to synchronize state information between PF firewalls. Its interfaces are assigned to physical interfaces with `ifconfig`. It is strongly recommended to set up pfsync on a separate network (or even VLAN), even if it is technically possible to lump in pfsync traffic with other traffic on a regular interface. The main reason for this recommendation is that pfsync itself does not do any authentication on its synchronization partners, and you can guarantee correct synchronization only if you are using dedicated interfaces for your pfsync traffic.

## ***The Project Specification: A Redundant Pair of Gateways***

To illustrate a useful failover setup with CARP and pfsync, we'll walk through an example, beginning with a network that currently has one gateway to the world. The following are our goals for the reconfigured network:

- The network work should keep functioning much the same way as earlier.
- We should have better availability with no noticeable downtime.
- The network should experience graceful failover with no interruption of active connections.

---

<sup>3</sup>. VRRP is described in RFC 2281 and RFC 3768. The patents involved are held by Cisco, IBM, and Nokia. See the RFCs for details.

We start with the relatively simple network from Chapter 3, which looks something like Figure 7-2.

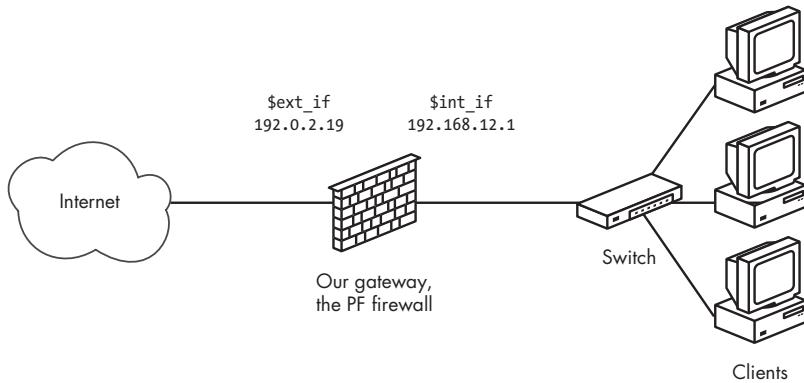


Figure 7-2: Network with a single gateway

We replace the single gateway with a redundant pair of gateways that share a private network between them for state-information updates over pfsync. The result is something like Figure 7-3.

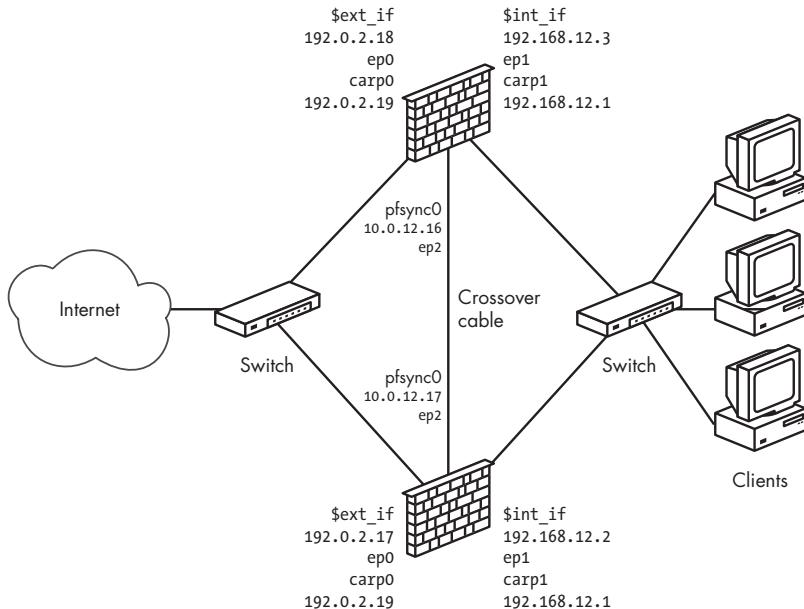


Figure 7-3: Network with redundant gateways

Note that the CARP addresses are virtual addresses. Unless you have console access to all machines in your CARP group, you will almost always want to assign an IP address to the physical interfaces. With a unique IP address for each of the physical interfaces, you will be able to communicate with the host and be absolutely sure with which machine you are interacting. Without IP addresses assigned to physical interfaces, you could find yourself with a setup

where the backup gateways are unable to communicate (except with hosts in networks where the physical interfaces have addresses assigned), until they become the master in the redundancy group and take over the virtual IP addresses.

By convention, the IP address assigned to the physical interface will belong in the same subnet as the virtual, shared IP address. By default, the kernel will try to assign the CARP address to a physical interface that is already configured with an address in the same subnet as the CARP address. You can make this interface selection explicit by specifying a different interface in the `carpdev` option in the `ifconfig` command string you use to set up the CARP interface.

**WARNING** *When you are reconfiguring your network and the default gateway address goes from being fixed to a specific interface and host to a virtual address, it's pretty much impossible to avoid temporary loss of connectivity.*

## **Setting Up CARP**

Most of the CARP setup lies in cabling (according to the schematic for your network), setting `sysctl` values, and issuing `ifconfig` commands. Also, on some systems, you will need to check that your kernel is set up with the required devices compiled in.

### **Checking Kernel Options**

On OpenBSD, both the CARP and pfsync devices are in the default `GENERIC` and `GENERIC.MP` kernel configurations. Unless you are running with a custom kernel where you removed these options, no kernel reconfiguration is necessary.

FreeBSD users need to check that your kernel has the CARP and pfsync devices compiled in. The `GENERIC` kernel does not contain these options by default. See the *FreeBSD Handbook* for information about how to compile and install a custom kernel with these options.

NetBSD users need to check that your kernel has pseudo-device CARP compiled in. NetBSD's default `GENERIC` kernel configuration does not have CARP compiled in. However, you will find the relevant line commented out in the `GENERIC` configuration file for easy inclusion. NetBSD does not yet support pfsync, due to claimed protocol-numbering issues that were unresolved at the time this chapter was written.

### **Setting sysctl Values**

On all CARP-capable systems, the basic functions are governed by a handful of `sysctl` variables. The main one, `net.inet.carp.allow`, is enabled by default. On a typical OpenBSD system, you will see this:

---

```
$ sysctl net.inet.carp.allow  
net.inet.carp.allow=1
```

---

This means that your system comes equipped for CARP.

If your kernel is not configured with a CARP device, this command will instead produce something like `sysctl: unknown oid 'net.inet.carp.allow'` on FreeBSD, or `sysctl: third level name 'carp' in 'net.inet.carp.allow' is invalid` on NetBSD.

Use this `sysctl` command to view all CARP-related variables:

---

```
$ sysctl net.inet.carp  
net.inet.carp.allow=1  
net.inet.carp.preempt=0  
net.inet.carp.log=2
```

---

**NOTE** *On FreeBSD, you will also encounter the variable `net.inet.carp.suppress_preempt`, which is a read-only status variable indicating whether or not preemption is possible. On systems with CARP code based on OpenBSD 4.2 or earlier, you will also see `net.inet.carp.arpbalance`, which is used to enable CARP ARP balancing to offer some limited load balancing for hosts on a local network.*

To enable the graceful failover between the gateways in the setup we are planning, we need to set the `net.inet.carp.preempt` variable:

---

```
$ sudo sysctl net.inet.carp.preempt=1
```

---

Setting the `net.inet.carp.preempt` variable means that on hosts with more than one network interface, such as our gateways, all CARP interfaces will move between master and backup status together. This setting needs to be identical on all hosts in the CARP group. When setting up, you need to repeat it on all hosts.

The `net.inet.carp.log` variable sets the debug level for CARP logging. The range of possible values is 0 through 7. The default is 2, which means only CARP state changes are logged.

### Setting Up Network Interfaces with ifconfig

Looking at the network diagram, we see that the local network uses addresses in the 192.168.12.0 network, while the external, Internet-facing interface is in the 192.0.2.0 network. With these address ranges and the CARP interface's default behavior in mind, the commands for setting up the virtual interfaces come naturally.

In addition to the usual network parameters, CARP interfaces require one extra parameter: the *virtual host ID* (`vhid`) that uniquely identifies the interfaces that are to share the virtual IP address.

**WARNING** *The `vhid` is an 8-bit value that must be set to a value that is unique within the network's broadcast domain. Setting the `vhid` to the wrong value can lead to several different types of network problems that can be hard to debug. There is even anecdotal evidence that redundancy and load-balancing systems based on VRRP use a virtual node identification scheme similar enough to CARP's that ID collisions with otherwise unrelated systems can occur and cause disruption.*

On the machine you want to set up as the initial master for the group, use these commands:

---

```
$ sudo ifconfig carp0 192.0.2.19 vhid 1
$ sudo ifconfig carp1 192.168.1.1 vhid 2
```

---

Note that we do not need to set the physical interface explicitly. The `carp0` and `carp1` virtual interfaces here will bind themselves to the physical interfaces that are already configured with addresses in the same subnets as the assigned CARP address.

**NOTE** *In contexts where the choice of physical network device for the CARP interface is not obvious from the existing network configuration, you can add a `carpdev` interface string to the `ifconfig` command. This can be useful in some nonintuitive configurations, as well as in scenarios where the number of free IP addresses in the relevant network is severely limited. At the time this chapter was written, the FreeBSD port of CARP did not offer the `carpdev` option.*

With `ifconfig`, you should be able to check that each CARP interface was properly configured:

---

```
$ ifconfig carp0
carp0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 00:00:5e:00:01:01
      carp: MASTER carpdev ep0 vhid 1 advbase 1 advskew 0
      groups: carp
      inet 192.0.2.19 netmask 0xffffffff broadcast 192.0.2.255
      inet6 fe80::200:5eff:fe00:101%carp0 prefixlen 64 scopeid 0x5
```

---

Note the `carp:` line, which indicates `MASTER` status.

On the backup, the setup is almost identical, except that you add the `advskew` parameter:

---

```
$ sudo ifconfig carp0 192.0.2.19 vhid 1 advskew 100
$ sudo ifconfig carp1 192.168.1.1 vhid 2 advskew 100
```

---

The `advskew` parameter indicates how much *less preferred* it is for the specified machine to take over than the current master. `advskew` and its companion value `advbase` are used to calculate the interval between the current host's announcements of its master status when it has taken over. The default value for `advbase` is 1; for `advskew`, the default value is 0. In our example, the master would announce every second ( $1 + 0/256$ ), while the backup would wait for  $1 + 100/256$  seconds. With `net.inet.carp.preempt=1`, when the master stops announcing or announces it is not available, the backups will take over, and the new master will start announcing at its configured rate.

Smaller `advskew` values mean shorter announcement intervals and higher likelihood for becoming the new master. If more hosts have the same `advskew`, the one that is already master will keep its master status.

From OpenBSD 4.1 onward, there is one more factor to the equation that determines which hosts takes over CARP master duty. The *demotion counter* is a value each CARP host announces for its interface group as a measure of readiness for its CARP interfaces. When the demotion counter value is 0, the host is in complete readiness; higher values indicate measures of degradation. You can set the demotion counter from the command line using `ifconfig -g`, but the value is usually set by the system itself, with higher values typically during the boot process. All other things equal, the host with the lowest demotion counter will win the contest to take over as the CARP master.

**NOTE** *At the time this chapter was written, the FreeBSD CARP port did not support setting the demotion counter.*

On the backup, use `ifconfig` once again to check that each CARP interface is properly configured:

---

```
$ ifconfig carp0
carp0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
    lladdr 00:00:5e:00:01:01
    carp: BACKUP carpdev ep0 vhid 1 advbase 1 advskew 100
    groups: carp
    inet 192.0.2.19 netmask 0xffffffff broadcast 192.0.2.255
    inet6 fe80::200:5eff:fe00:101%carp0 prefixlen 64 scopeid 0x5
```

---

Here, the output is only slightly different. Notice that the `carp:` line indicates BACKUP status along with the `advbase` and `advskew` parameters.

For actual production use, you will want to add a measure of security against unauthorized CARP activity by configuring the members of the CARP group with a shared, secret passphrase, such as the following:<sup>4</sup>

---

```
$ sudo ifconfig carp0 pass mekmitasdigoat 192.0.2.19 vhid 1
$ sudo ifconfig carp1 pass mekmitasdigoat 192.168.1.1 vhid 2
```

---

**WARNING** *Just like any other password, the passphrase will then be a required ingredient in all CARP traffic in your setup. Take care to configure all CARP interfaces in a failover group with the same passphrase (or none).*

After you have figured out the appropriate settings, you will want to preserve them through future system reboots by putting them in the proper files in `/etc`:

- On OpenBSD, put the proper `ifconfig` parameters into `hostname.carp0` and `hostname.carp1`.
- On FreeBSD and NetBSD, put the relevant lines in your `rc.conf` file as contents of the `ifconfig_carp0=` and `ifconfig_carp1=` variables.

---

4. This particular passphrase has a very specific meaning. A web search will reveal its significance and why it is *de rigueur* for modern networking documentation. The definitive answer can be found via the `openbsd-misc` mailing list archives.

## **Keeping States Synchronized: Adding pfsync**

The final piece of the configuration is to set up state-table synchronization between the hosts in your redundancy group. With synchronized state tables on the redundant firewalls, in almost all cases, there will be no apparent traffic disruption during failover. This is accomplished through a set of properly configured pfsync interfaces. (As noted earlier, at the time of writing, NetBSD does not support pfsync.)

Configuring pfsync interfaces is a matter of some planning and a few fairly straightforward ifconfig commands. It is possible to set up pfsync on any configured network interface, but it is generally a better idea to set up a separate network for the synchronization.

In our sample configuration (Figure 7-3), we have set aside a tiny network for the purpose. Here, a crossover cable connects the two Ethernet interfaces, but in configurations with more than two hosts in the failover group, you may want to set up with a separate switch, hub, or VLAN.

In our sample configuration, the interfaces we are planning to use for the synchronization have already been assigned the IP addresses 10.0.12.16 and 10.0.12.17, respectively. With the basic TCP/IP configuration in place, the complete pfsync setup for each of the two synchronization partner interfaces is as follows:

---

```
$ sudo ifconfig pfsync0 syncdev ep2
```

---

This illustrates the advantage of identical hardware configurations, as well as keeping pfsync traffic on a physically separate network.

The pfsync protocol itself offers little in the way of security features. It has no authentication mechanism and, by default, communicates via IP multicast traffic. However, for the cases where a physically separate network is not a feasible option, you can tighten up your pfsync security by setting up pfsync to synchronize only with a specified syncpeer:

---

```
$ sudo ifconfig pfsync0 syncpeer 10.0.12.16 syncdev ep2
```

---

This would produce a configured interface that shows up in your ifconfig output like this:

---

```
pfsync0: flags=41<UP,RUNNING> mtu 1500
        priority: 0
        pfsync: syncdev: ep2 syncpeer: 10.0.12.16 maxupd: 128 defer: off
        groups: carp pfsync
```

---

One other option would be to set up an IPsec tunnel and use that to protect the sync traffic. The ifconfig command would then be as follows:

---

```
$ sudo ifconfig pfsync0 syncpeer 10.0.12.16 syncdev enc0
```

---

This means that the syncdev device becomes the enc0 encapsulating interface instead of the physical interface.

**NOTE** *If possible, set up your synchronization to happen across a physically separate, dedicated network. Any lost pfsync updates could lead to a less than clean failover.*

A very useful way to check that your PF state synchronization is running properly is to watch the state table on both (or all) synchronized hosts using **sysat states** on each machine. The command gives you a live display of states, where you can see updates happening in bulk on the sync targets. In between the synchronizations, states should display identically on all hosts, except that traffic counters (such as the number of packets and bytes passed) will display updates only on the host that handles the actual connection.

This takes us to the end of the basic network configuration for CARP-based failover. In the next section, we consider what you need to keep in mind when writing rule sets for redundant configurations.

### ***Putting Together a Rule Set***

After all the contortions we've been through in order to get the basic networking configured, you are probably wondering what it will take to migrate the rules you have put in your current *pf.conf* to the new setup. You'll be glad to know that it won't take very much. The main change you have introduced is essentially invisible to the rest of the world, and a well-designed rule set for a single gateway configuration will generally work well for a redundant setup, too.

However, you have introduced two additional protocols: CARP and pfsync. In all likelihood, you will need to make some relatively minor changes to your rule set in order to let the failover work properly. Basically, you need to pass the CARP and pfsync traffic to the appropriate interfaces.

The most readable way to handle the CARP traffic is to introduce a macro definition for your *carpdevs* and an accompanying *pass* rule, such as the following:

---

```
pass on $carpdevs proto carp
```

---

You need to pass pfsync traffic on the appropriate interfaces.

Similary, for pfsync traffic, you can introduce a macro definition for your *syncdev* and an accompanying *pass* rule:

---

```
pass on $syncdev proto pfsync
```

---

If you want to take the pfsync device out of the filtering equation altogether, use this rule:

---

```
set skip on $syncdev
```

---

Skipping the pfsync interfaces entirely for filtering is cheaper performance-wise than filtering and passing.

Also, you should consider the roles of the virtual CARP interface and its address versus the physical interface. As far as PF is concerned, all traffic will pass through the physical interfaces, but the traffic may have the CARP interface's IP addresses as source or destination addresses.

You may have rules in your configuration that you don't want to bother to synchronize in case of a failover, such as connections to services that run on the gateway itself. One prime example would be the typical rule to allow SSH in for the administrator:

---

```
pass in on $int_if from $ssh_allowed to self
```

---

For those rules, you could use the state option `no-sync` to prevent synchronizing state changes for connections that are really not relevant after the failover has happened:

---

```
pass in on $int_if from $ssh_allowed to self keep state (no-sync)
```

---

With this configuration in hand, you will be able to schedule operating system upgrades and similar former downtime-producing activities on members of your CARPed group of systems at times when they are convenient to the system administrator, with no measurable or noticeable downtime for your services' users.

## IFSTATED, THE INTERFACE STATE DAEMON

In properly CARPed setups, the basic networking functionality is very well provided for. However, outside the basic networking, your setup may include elements that need special attention when the network configuration on a host changes. For example, some services might need to be started or stopped when a specific interface goes down or starts up again, or you may want to run specific commands or scripts in response to interface state changes. If this sounds familiar, ifstated is for you.

ifstated was introduced in OpenBSD 3.5 as a tool to trigger actions based on changes in network interfaces' state. You can think of ifstated as "CARP's little helper."

ifstated is in the base system on OpenBSD and available via ports as `net/ifstated` on FreeBSD. On OpenBSD, `/etc/ifstated.conf` (or `/usr/local/etc/ifstated.conf` if you installed the port on FreeBSD) contains an almost-ready-to-run configuration that gives you more than a few pointers on how to set up ifstated for a CARPed environment.

The main controlling objects are interfaces and their states, where `carp0.link.up` is the state where the `carp0` interface has a link, and you perform actions in response to state changes. These are specified in a straightforward scripting language with basic features like variables, macros, and simple logical conditionals. See `man ifstated` and `man ifstated.conf` if you want to dig deeper into the topic.

## **CARP for Load Balancing**

Redundancy by failover is nice, but in some contexts, it is less attractive to have hardware just sitting around in case of failure, and you may prefer to create a configuration that spreads the network load over several hosts.

In addition to ARP balancing, which works by calculating hashes based on the source MAC address on incoming connections, CARP in OpenBSD 4.3 onward supports several varieties of IP-based load balancing, where traffic is allocated based on hashes calculated from the connections' source and destination IP addresses. Since ARP balancing is based on the source MAC address, it will work only for hosts in the directly connected network segment. The IP-based methods are appropriate for load-balancing connections to and from the Internet at large.

Which mode is the wise choice for your application depends on the exact specifications of the rest of the network equipment you need to work with. The basic ip balancing mode uses a multicast MAC address to make the directly connected switch forward traffic to all hosts in the load-balancing cluster. The combination of a unicast IP address and a multicast MAC address is not supported by some systems. In those cases, you may need to configure your load balancing in ip-unicast mode (which uses a unicast MAC address) and configure your switch to forward to the appropriate hosts, or in ip-stealth mode, which does not use the multicast MAC address at all. As usual, the devil is in the details, and the answers are found in man pages and other documentation, most likely with a bit of experimentation thrown in.

**NOTE** *Traditionally, relayd has been used to do intelligent load balancing as the front end for servers that offer services to the rest of the world. In OpenBSD 4.7, relayd acquired the ability to keep track of available uplinks and alter the system's routing tables based on link health, with the functionality wrapped in a bundle with the router keyword. For setups with several possible uplinks or variant routing tables, you can set up relayd to choose your uplink or, with a little help from the sysctl variables net.inet.ip.multipath and net.inet6.ip6.multipath, perform load balancing across available routes and uplinks. The specifics will vary with your networking environment. The relayd.conf man page contains a complete example to get you started.*

In load-balancing mode, the CARP concept is extended by letting each CARP interface be a member of not only one failover group, but as many load-balancing groups as there are physical hosts that will share the virtual address. In contrast with the failover case where there can be only one master, each node in a load balancing cluster *must* be the master of its own group so it can receive traffic. Which group, and by extension which physical host, ends up handling a given connection is determined by CARP via a hash value calculation based on the connection's source MAC address in the ARP balancing case, and source and destination IP address in the IP balancing case, in addition to actual availability.

The downside of this is that each of these groups consumes one virtual host ID, so you will run out of these IDs quite a bit quicker in a load-balancing configuration than with failover only. In fact, there is a hard upper limit to the number of CARP-based load balancing *clusters* at 32 virtual host IDs.

The `advskev` parameter plays a similar role in load-balancing configurations as in the failover ones, but the `ifconfig` (and `hostname.carpN`) syntax for CARP load balancing is slightly different from the failover case.

Changing the CARP failover group we built up over the previous sections to a load-balancing cluster is as easy as editing the configuration files and reloading. In this example, we go for an IP load-balancing scheme. If you choose a different load-balancing scheme, the configuration itself differs only in the keyword for mode selection.

On the first host, we change `/etc/hostname.carp0` to this:

---

```
pass mekmitasdigoat 192.0.2.19 balancing ip carpnodes 5:100,6:0
```

---

This means that this on this host, the `carp0` interface is a member of the group with `vhid 5`, where it has an `advskev` of `100` as well as the one with `vhid 6`, where it is the prime candidate for becoming initial master, with an `advskev` set to `0`.

Next, we change `/etc/hostname.carp1` to read as follows:

---

```
pass mekmitasdigoat 192.168.12.1 balancing ip carpnodes 3:100,4:0
```

---

For `carp1`, the memberships are `vhids 3` and `4`, with `advskev` values of `100` and `0`, respectively.

For the other host, the `advskev` values are reversed, but the configuration is otherwise predictably similar. Here, `/etc/hostname.carp0` reads as follows:

---

```
pass mekmitasdigoat 192.0.2.19 balancing ip carpnodes 5:0,6:100
```

---

It is a member of `vhid 5` with `advskev 0`, and a member of `vhid 6` with `advskev 100`. Complementing this is the `/etc/hostname.carp1` file that reads like this:

---

```
pass mekmitasdigoat 192.168.12.1 balancing ip carpnodes 3:0,4:100
```

---

Again, `carp1` is a member of `vhid 3` and `4`, with `advskev 0` in the first and `100` in the other.

The `ifconfig` output for the `carp` interface group on the first host looks like this:

---

```
$ ifconfig carp
carp0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 01:00:5e:00:01:05
      priority: 0
      carp: carpdev vr0 advbase 1 balancing ip
            state MASTER vhid 5 advskev 0
            state BACKUP vhid 6 advskev 100
      groups: carp
      inet 192.0.2.19 netmask 0xffffffff broadcast 192.0.2.255
      inet6 fe80::200:24ff:fe:1c10%carp0 prefixlen 64 scopeid 0x7
```

---

```
carp1: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 01:00:5e:00:01:03
      priority: 0
      carp: carpdev vr1 advbase 1 balancing ip
            state MASTER vhid 3 advskew 0
            state BACKUP vhid 4 advskew 100
      groups: carp
      inet 192.168.12.1 netmask 0xffffffff broadcast 192.168.12.255
      inet6 fe80::200:24ff:fe:1c10%carp1 prefixlen 64 scopeid 0x8
pfsync0: flags=41<UP,RUNNING> mtu 1500
      priority: 0
      pfsync: syncdev: vr2 syncpeer: 10.0.12.17 maxupd: 128 defer: off
      groups: carp pfsync
```

---

The other host has the following ifconfig output:

```
$ ifconfig carp
carp0: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 01:00:5e:00:01:05
      priority: 0
      carp: carpdev vr0 advbase 1 balancing ip
            state BACKUP vhid 5 advskew 100
            state MASTER vhid 6 advskew 0
      groups: carp
      inet 192.0.2.19 netmask 0xffffffff broadcast 192.0.2.255
      inet6 fe80::200:24ff:fe:1c18%carp0 prefixlen 64 scopeid 0x7
carp1: flags=8843<UP,BROADCAST,RUNNING,SIMPLEX,MULTICAST> mtu 1500
      lladdr 01:00:5e:00:01:03
      priority: 0
      carp: carpdev vr1 advbase 1 balancing ip
            state BACKUP vhid 3 advskew 100
            state MASTER vhid 4 advskew 0
      groups: carp
      inet 192.168.12.1 netmask 0xffffffff broadcast 192.168.12.255
      inet6 fe80::200:24ff:fe:1c18%carp1 prefixlen 64 scopeid 0x8
pfsync0: flags=41<UP,RUNNING> mtu 1500
      priority: 0
      pfsync: syncdev: vr2 syncpeer: 10.0.12.16 maxupd: 128 defer: off
      groups: carp pfsync
```

---

If we had three nodes in our load-balancing scheme, each carp interface would need to be a member of an additional group, for a total of three groups. In short, for each physical host you introduce in the load-balancing group, each carp interface becomes the member of an additional group.

Once you have set up the load-balancing cluster, you can check the actual flow of connections by running **systat states** on each of the hosts in your load-balancing cluster. I recommend doing just that for a few minutes to make sure that the system works as expected and to see that all the effort you put in has been worth it.

# 8

## **LOGGING, MONITORING, AND STATISTICS**



Exercising control over a network—whether for your home networking needs or in a professional context—is likely to be a main objective for anyone who reads this book.

One necessary element of keeping control is having access to all relevant information about what happens in your network. Fortunately for us, PF (like most components of Unix-like systems) is able to generate log data for network activity.

PF offers a wealth of options for setting the logging detail level, processing log files, and extracting specific kinds of data. You can already do a lot with the tools that are in your base system, and several other tools are available via your package system to collect, study, and view log data in a number of useful ways. In this chapter, we take a closer look at PF logs in general and some of the tools you can use to extract and present information.

## PF Logs: The Basics

The information that PF logs and the level of logging detail is up to you, as determined by your rule set. Basic logging is simple: For each rule that you want log data for, add the `log` keyword. When you load the rule set with `log` added to one or more rules, any packet that starts a connection matching the rule (blocked, passed, or matched) is copied to a `pflog` device. PF will also store certain additional data such as the timestamp, interface, whether the packet was blocked or passed, and the associated rule number from the loaded rule set.

PF log data is collected by the `pflogd` logging daemon, which is started by default when PF is enabled at system startup. The default location for storing the log data is `/var/log/pflog`. The log is written in a binary format intended to be read by `tcpdump`. We'll discuss additional tools to extract and display information from your log file later. The log file format is a well documented and widely supported binary format.

To get started, here's a basic log example. Start with the rules you want to log and add the `log` keyword:

---

```
block log
pass log quick proto { tcp, udp } to port ssh
```

---

Reload the rule set, and you should see the timestamp on your `/var/log/pflog` file change as the file starts growing. To see what is being stored there, use `tcpdump` with the `-r` option to read the file.

If logging has been going on for a while, entering the following on a command line can produce large amounts of output.

---

```
$ sudo tcpdump -n -ttt -r /var/log/pflog
```

---

For example, the following are just the first lines from a file several screens long, with almost all lines long enough to wrap:

---

```
$ sudo tcpdump -n -ttt -r /var/log/pflog
tcpdump: WARNING: snaplen raised from 96 to 116
Sep 13 13:00:30.5556038 rule 10/(match) pass in on epic0: 194.54.107.19.34834 > 194.54.103.66.113: S
3097635127:3097635127(0) win 16384 <mss 1460,nop,nop,sackOK,nop,wscale 0,[|tcp]> (DF)
Sep 13 13:00:30.5556063 rule 10/(match) pass out on fxp0: 194.54.107.19.34834 > 194.54.103.66.113: S
3097635127:3097635127(0) win 16384 <mss 1460,nop,nop,sackOK,nop,wscale 0,[|tcp]> (DF)
Sep 13 13:01:07.796096 rule 10/(match) pass in on epic0: 194.54.107.19.29572 > 194.54.103.66.113: S
2345499144:2345499144(0) win 16384 <mss 1460,nop,nop,sackOK,nop,wscale 0,[|tcp]> (DF)
Sep 13 13:01:07.796120 rule 10/(match) pass out on fxp0: 194.54.107.19.29572 > 194.54.103.66.113: S
2345499144:2345499144(0) win 16384 <mss 1460,nop,nop,sackOK,nop,wscale 0,[|tcp]> (DF)
Sep 13 13:01:15.096643 rule 10/(match) pass in on epic0: 194.54.107.19.29774 > 194.54.103.65.53:
49442 [1au][|domain]
Sep 13 13:01:15.607619 rule 12/(match) pass in on epic0: 194.54.107.19.29774 > 194.54.107.18.53:
34932 [1au][|domain]
```

---

The `tcpdump` program is very flexible, especially when it comes to output, and it offers a number of display choices. The format in this example follows from the options we fed to `tcpdump`. The program almost always displays the date and time the packet arrived (the `-ttt` option specifies this long format). Next, `tcpdump` lists the rule number in the loaded rule set, the interface on which the packet appeared, the source and target address and ports (the `-n` option tells `tcpdump` to display IP addresses, not hostnames), and the various packet properties.

**NOTE** *The rule numbers in your log files refer to the loaded, in-memory rule set. Your rule set goes through some automatic steps during the loading process, such as macro expansion and optimizations, which make it likely that the rule number as stored in the logs will not quite match what you would find by counting from the top of your pf.conf file. If it isn't immediately obvious to you which rule matched, use `pfctl -vvv rules` and study the output.*

In our `tcpdump` output example, we see that the tenth rule (rule 10) in the loaded rule set seems to be a catchall that matches both IDENT requests and domain name lookups. This is the kind of output you will find invaluable when debugging, and it is essential to have this kind of data available in order to stay on top of your network. With a little effort and careful reading of the `tcpdump` man pages, you should be able to extract useful information from your log data.

For a live display of the traffic you log, use `tcpdump` to read log information directly from the log device. To do so, use the `-i` option to specify which interface you want `tcpdump` to read from, as follows:

---

```
$ sudo tcpdump -netti pflog0
tcpdump: WARNING: pflog0: no IPv4 address assigned
tcpdump: listening on pflog0, link-type PFLOG
Sep 13 15:26:52.122002 rule 17/(match) pass in on epic0: 91.143.126.48.46618 >
194.54.103.65.22: [|tcp] (DF)
Sep 13 15:28:02.771442 rule 12/(match) pass in on epic0: 194.54.107.19.8025 >
194.54.107.18.8025: udp 50
Sep 13 15:28:02.773958 rule 10/(match) pass in on epic0: 194.54.107.19.8025 >
194.54.103.65.8025: udp 50
Sep 13 15:29:27.882888 rule 10/(match) pass in on epic0: 194.54.107.19.29774 >
194.54.103.65.53:[|domain]
Sep 13 15:29:28.394320 rule 12/(match) pass in on epic0: 194.54.107.19.29774 >
194.54.107.18.53:[|domain]
```

---

This sequence begins with an SSH connection. The next two connections are `spamd` synchronizations, followed by two domain name lookups. If you were to leave this command running, the displayed lines would eventually scroll off the top of your screen, but you could redirect the data to a file or to a separate program for further processing.

**NOTE** *Sometimes you will be interested mainly in traffic between specific hosts, or traffic matching specific criteria. For these cases, `tcpdump`'s filtering features can be useful. See `man tcpdump` for details.*

## **Logging All Packets: log (all)**

For most debugging and lightweight monitoring purposes, logging the first packet in a connection provides enough information. However, sometimes you may want to log all packets that match certain rules. To do so, use the (all) logging option in the rules you want to monitor. After making this change to our minimal rule set, we have this:

---

```
block log (all)
pass log (all) quick proto tcp to port ssh keep state
```

---

This option makes the logs quite a bit more verbose. To illustrate just how much more data log (all) generates, we'll use the following rule set fragment, which passes domain name lookups and network time synchronizations:

---

```
udp_services = "{ domain, ntp }"
pass log (all) inet proto udp to port $udp_services
```

---

With these rules in place, here's an example of what happens when a Russian name server sends a domain name request to a server in our network:

---

```
$ sudo tcpdump -n -ttt -i pflog0 port domain
tcpdump: WARNING: pflog0: no IPv4 address assigned
tcpdump: listening on pflog0, link-type PFLOG
Sep 30 14:27:41.260190 212.5.66.14.53 > 194.54.107.19.53:[|domain]
Sep 30 14:27:41.260253 212.5.66.14.53 > 194.54.107.19.53:[|domain]
Sep 30 14:27:41.260267 212.5.66.14.53 > 194.54.107.19.53:[|domain]
Sep 30 14:27:41.260638 194.54.107.19.53 > 212.5.66.14.53:[|domain]
Sep 30 14:27:41.260798 194.54.107.19.53 > 212.5.66.14.53:[|domain]
Sep 30 14:27:41.260923 194.54.107.19.53 > 212.5.66.14.53:[|domain]
```

---

We now have six entries instead of just one.

### **LOG RESPONSIBLY!**

Creating logs of any kind could have surprising consequences, including some legal implications. Once you start storing log data generated by your network traffic, you are creating a store of information about your users. There may be good technical and business reasons to store logs for extended periods, but logging just enough data for sufficient time can be a fine art.

You probably have some idea of the practical issues related to generating log data, such as arranging for sufficient storage to retain enough log data long enough to be useful. The legal implications will vary according to your location. Some countries and territories have specific requirements for handling log data, along with restrictions on how that data may be used and how long logs can be retained. Others require service providers to retain traffic logs for a specific period of time, in some cases with a requirement to deliver any such data to law enforcement upon request. Make sure you understand the legal issues before you build a logging infrastructure.

Even with all but port domain filtered out by tcpdump, adding log (all) to one or more rules considerably increases the amount of data in your logs. If you need to log all traffic, but your gateway's storage capacity is limited, you may find yourself shopping for additional storage. Also, recording and storing traffic logs with this level of detail is likely to have legal implications.

## ***Logging to Several pflog Interfaces***

PF versions older than OpenBSD 4.1 offered only one pflog interface. That changed with OpenBSD 4.1, when the pflog interface became a *cloneable* device, meaning that you can use ifconfig commands to create several pflog interfaces, in addition to the default pflog0. This makes it possible to record the log data for different parts of your rule set to separate pflog interfaces, and easier to process the resulting data separately if necessary.

The required changes to your setup are subtle but effective. To log to several interfaces, make sure that all the log interfaces your rule set uses are created. You don't need to create the devices before the rule set is loaded; if your rule set logs to a nonexistent interface, the log data is simply discarded.

When tuning your setup to use several pflog interfaces, you will most likely add the required interfaces from the command line, like so:

---

```
$ sudo ifconfig create pflog1
```

---

Specify the log device when you add the log keyword to your rule set, as follows:

---

```
pass log (to pflog1) proto tcp to $emailserver port $email  
pass log (to pflog1) proto tcp from $emailserver to port smtp
```

---

For a more permanent configuration on OpenBSD, create a *hostname.pflog1* file containing only up, and similar *hostname.pflogN* files for any additional logging interfaces.

On FreeBSD, the configuration of the cloned pflog interfaces belongs in your *rc.conf* file, in the following form:

---

```
ifconfig_pflog1="up"
```

---

On NetBSD, cloning pflog interfaces is not an option as of this writing.

As you saw in Chapter 6, directing log information for different parts of your rule set to separate interfaces makes it possible to feed different parts of the log data PF produces to separate applications. This makes it easier to have programs like spamlogd process only the relevant information, while you feed other parts of your PF log data to other log-processing programs.

## ***Logging to Syslog, Local or Remote***

One way to avoid storing PF log data on the gateway itself is to instruct your gateway to log to another machine. If you already have a centralized logging infrastructure in place, this is a fairly logical thing to do, even if

PF's ordinary logging mechanisms were not really designed with traditional syslog-style logging in mind.

As any old BSD hand will tell you, the traditional syslog system log facility is a bit naïve about managing the data it receives over UDP from other hosts, with denial-of-service attacks involving full disks one frequently mentioned danger. There is also the ever-present risk that log information will be lost under high load on either individual systems or the network. Therefore, consider setting up remote logging *only* if all hosts involved communicate over a well-secured network. On most BSDs, syslogd is not set up by default to accept log data from other hosts. (See the `syslogd` man page for information about how to enable listening for log data from remote hosts if you plan to use remote syslog logging.)

If you would still like to do your PF logging via syslog, the following is a short recipe for how to accomplish this. In ordinary PF setups, `pflogd` copies the log data to the log file. When you primarily want to store the log data on a remote system, you should disable `pflogd`'s data accumulation by changing its log file to `/dev/null`, via the daemon's startup options in `rc.conf.local` (on OpenBSD), like so:

---

```
pflogd_flags="-f /dev/null"
```

---

On FreeBSD and NetBSD, change the `pflog_logfile=` line in `rc.conf` as follows, and then kill and restart the `pflogd` process with its new parameters:

---

```
pflog_logfile="/dev/null"
```

---

Next, make sure that the log data, now no longer collected by `pflogd`, is transmitted in a meaningful way to your log-processing system instead. This step has two parts: First, set up your system logger to transmit data to the log processing system, and then use `tcpdump` with `logger` to convert the data and inject it into the syslog system.

To set up `syslogd` to process the data, choose your *log facility*, *log level*, and *action* and put the resulting line in `/etc/syslog.conf`. These concepts are very well explained in `man syslog.conf`, which is required reading if you want to understand system logs. The *action* part is usually a file in a local filesystem. For example, if you have already set up the system logger at `loghost.example.com` to receive your data, choosing log facility `local2` with log level `info`, enter this line:

---

local2.info	@loghost.example.com
-------------	----------------------

---

Once you've made this change, restart `syslogd` to make it read the new settings.

Next, set `tcpdump` to convert the log data from the `pflog` device and feed it to `logger`, which will then send it to the system logger. Here, we reuse the `tcpdump` command from the basic examples earlier in this chapter, with some useful additions:

---

```
$ sudo nohup tcpdump -l nettti pflog0 | logger -t pf -p local2.info &
```

---

The `nohup` command makes sure the process keeps running even if it does not have a controlling terminal or it's put in the background (as we do here with the trailing &). The `-l` option to the `tcpdump` command specifies line-buffered output, which is useful for redirecting to other programs. The `logger` option adds the tag `pf` to identify the PF data in the stream and specifies log priority with the `-p` option as `local2.info`. The result is logged to the file you specify on the logging host, with entries that will look something like this:

---

```
pf: Sep 21 14:05:11.492590 rule 93/(match) pass in on ath0: 10.168.103.11.15842 >
82.117.50.17.80: [|tcp] (DF)
pf: Sep 21 14:05:11.492648 rule 93/(match) pass out on xlo: 194.54.107.19.15842 >
82.117.50.17.80: [|tcp] (DF)
pf: Sep 21 14:05:11.506289 rule 93/(match) pass in on ath0: 10.168.103.11.27984 >
82.117.50.17.80: [|tcp] (DF)
pf: Sep 21 14:05:11.506330 rule 93/(match) pass out on xlo: 194.54.107.19.27984 >
82.117.50.17.80: [|tcp] (DF)
pf: Sep 21 14:05:11.573561 rule 136/(match) pass in on ath0: 10.168.103.11.6430 >
10.168.103.1.53:[|domain]
pf: Sep 21 14:05:11.574276 rule 136/(match) pass out on xlo: 194.54.107.19.26281 >
209.62.178.21.53:[|domain]
```

---

This log fragment shows mainly web-browsing activities from a client in a NATed local network as seen from the gateway's perspective, with accompanying domain name lookups.

### ***Tracking Statistics for Each Rule with Labels***

The sequential information you get from retrieving log data basically tracks packet movements over time. In other contexts, the sequence or history of connections is less important than aggregates, such as number of packets or bytes that have matched a rule since the counters were last cleared.

At the end of Chapter 2, you saw how to use `pfctl -s info` to view the global aggregate counters, along with other data. For a more detailed breakdown of the data, track traffic totals on a per-rule basis with a slightly different form of `pfctl` command, such as `pfctl -vs rules`, to display statistics along with the rule, as shown here:

---

```
$ pfctl -vs rules
pass inet proto tcp from any to 192.0.2.225 port = smtp flags S/SA keep state label "mail-in"
[ Evaluations: 1664158    Packets: 1601986    Bytes: 763762591    States: 0      ]
[ Inserted: uid 0 pid 24490 ]
pass inet proto tcp from 192.0.2.225 to any port = smtp flags S/SA keep state label "mail-out"
[ Evaluations: 2814933    Packets: 2711211    Bytes: 492510664    States: 0      ]
[ Inserted: uid 0 pid 24490 ]
```

---

The format of this output is easy to read, and obviously designed for contexts where you want to get an idea of what is going on at a glance. If you specify even more verbose output with `pfctl -vvs rules`, you will see essentially the same display, with rule numbers added. On the other hand, the output from this command is not very well suited for feeding to a script or other program for further processing. To extract these statistics and a few more items in

a slightly more script-friendly format (and make your own decisions about which rules are worth tracking), use rule *labels*.

Labels do more than identify rules for processing specific kinds of traffic; they also make it easier to extract the traffic statistics. By attaching labels to rules, you can store certain extra data about parts of your rule set. For example, you could use labeling to measure bandwidth use for accounting purposes.

In the following example, we attach the labels `mail-in` and `mail-out` to our `pass` rules for incoming and outgoing mail traffic, respectively.

---

```
pass log proto { tcp, udp } to $emailserver port smtp label "mail-in"
pass log proto { tcp, udp } from $emailserver to port smtp label "mail-out"
```

---

Once you have loaded the rule set with labels, check the data using `pfctl -vsl`:

---

```
$ sudo pfctl -vsl
 ❶   ❷   ❸   ❹   ❺   ❻   ❷   ❸
mail-in 1664158 1601986 763762591 887895 682427415 714091 81335176
mail-out 2814933 2711211 492510664 1407278 239776267 1303933 252734397
```

---

This output contains the following information:

- ❶ The label
- ❷ The number of times the rule has been evaluated
- ❸ The total number of packets passed
- ❹ The total number of bytes passed
- ❺ The number of packets passed in
- ❻ The number of bytes passed in
- ❷ The number of packets passed out
- ❸ The number of bytes passed out

The format of this list makes it very well suited for parsing by scripts and applications.

The labels accumulate data from the time the rule set is loaded until their counters are reset. And, in many contexts, it makes sense to set up a cron job that reads label values at fixed intervals, and then puts those values into permanent storage.

If you choose to run the data collection at fixed intervals, consider collecting the data using `pfctl -vszl`. The `z` option resets the counters once `pfctl` has read them, with the result that your data collector will then be fetching *periodic data*, accumulated since the command or the script was last run.

**NOTE** *Rules with macros and lists expand to several distinct rules. If your rule set contains rules with lists and macros that have a label attached, the in-memory result will be a number of rules, each with a separate, identically named label attached to it. While this*

*may lead to confusing sudo pfctl -vsl output, it shouldn't be a problem as long as the application or script that receives the data can interpret the data correctly by adding up the totals for the identical labels.*

## Additional Tools for PF Logs and Statistics

One other important component of staying in control of your network is having the ability to keep an updated view of your system's status. In this section, we'll examine a selection of monitoring tools that you may find useful. All the tools presented here are available either in the base system or via the package system on OpenBSD and FreeBSD (and with one exception, on NetBSD).

### ***Keeping an Eye on Things with systat***

If you are interested in seeing an instant snapshot of the traffic passing through your systems right now, the `systat` program on OpenBSD offers several useful views. In Chapter 7, we looked briefly at `systat` queues to see how traffic was assigned to queues in our ALTQ rule set. Here, we will review some additional useful options.

The `systat` program is available on all BSD operating systems, in slightly different versions. On all systems, `systat` offers views of system statistics, with some minor variations in syntax and output. For example, the `queues` view is one of several `systat` views available in recent OpenBSD versions, but not in FreeBSD or NetBSD as of this writing.

For a more general view of the current state table than that offered by `queues`, try `systat states`, which gives a listing very similar to the `top(1)` process listing. Here is an example of typical `systat states` output:

2 users		Load 0.24 0.28 0.27 (1-16 of 895)			Fri Apr 1 14:00:04 2011									
PR	D	SRC	DEST	STATE	AGE	EXP	PKTS	BYTES	RATE	PEAK	Avg	RU	G	
udp	O	192.168.103.1:56729	192.168.103.9:12345	1:0	8340m	25	372K	542M	1492	4774	1137	*		
tcp	I	10.168.103.15:47185	213.187.179.198:22	4:4	62377	86398	2954	613K	13264	23654	10	18		
tcp	I	10.168.103.15:2796	213.187.179.198:22	4:4	62368	86219	4014	679K	0	0	11	18		
tcp	I	10.168.103.15:15599	129.240.64.10:6667	4:4	61998	86375	9266	849K	0	58	14	*		
tcp	O	213.187.179.198:1559	129.240.64.10:6667	4:4	61998	86375	9266	849K	0	58	14	*	1	
tcp	I	10.168.103.15:8923	140.211.166.4:6667	4:4	61843	86385	15677	4794K	0	299	79	*		
tcp	O	213.187.179.198:8923	140.211.166.4:6667	4:4	61843	86385	15677	4794K	0	299	79	*	1	
tcp	I	10.168.103.15:47047	217.17.33.10:6667	4:4	61808	86385	7093	556K	0	88	9	*		
tcp	O	213.187.179.198:4704	217.17.33.10:6667	4:4	61808	86385	7093	556K	0	88	9	*	1	
tcp	I	10.168.103.15:30006	203.27.221.42:6667	4:4	61744	86375	6000	487K	0	49	8	*		
tcp	O	213.187.179.198:3000	203.27.221.42:6667	4:4	61744	86375	6000	487K	0	49	8	*	1	
tcp	I	10.168.103.15:31709	209.250.145.51:6667	4:4	61744	86385	6646	613K	0	114	10	*		
tcp	O	213.187.179.198:3170	209.250.145.51:6667	4:4	61744	86385	6646	613K	0	114	10	*	1	
tcp	I	192.168.103.254:5386	69.90.74.197:80	4:4	56718	29844	10	3282	0	0	0	*		
tcp	O	213.187.179.198:5386	69.90.74.197:80	4:4	56718	29844	10	3282	0	0	0	*	1	
tcp	I	10.168.103.15:33241	192.168.103.84:22	4:4	46916	82678	7555	897K	0	0	19	*		

If your states do not fit on one screen, just page through the live display.

Similarly, `systat rules` displays a live view of packets, bytes, and other statistics for your loaded rule set, as in this example:

---

2 users	Load 1.25 0.87 0.52 (1-16 of 239)	Fri Apr 1 14:01:59 2011
<hr/>		
RUL ANCHOR	A DIR L Q IF PR K PKTS BYTES STATE MAX INFO	
0	M In 26M 12G 4946K all max-mss 1440	
1	M Out nfe0 4853K 3162M 94858 inet from 10.0.0.0/8 to any queue(q_def	
2	M Out nfe0 3318K 2430M 61672 inet from 192.168.103.0/24 to any queue	
3	M Out nfe0 tcp 6404K 4341M 134K from any to any port = www queue(q_web,	
4	M Out nfe0 tcp 84298 43M 1594 from any to any port = https queue(q_we	
5	M Out nfe0 tcp 502 34677 63 from any to any port = domain queue(q_d	
6	M Out nfe0 udp 512K 64M 257K from any to any port = domain queue(q_d	
7	M Out nfe0 icmp 11 1008 3 all queue(q_dns, q_pri)	
8	B Any L 14638 1346K 0 return all	
9	B Any Q 95 5628 0 return from <bruteforce> to any	
10	P Any 1139K 1005M 757 all flags any	
11	P In Q tcp K 18538 1350K 708 inet from any to any port = ftp	
12	P Out tcp K 0 0 0 inet from 127.0.0.1/32 to any port = ftp	
13	P Any 1421 128K 134 all flags any	
14	P In L egress tcp K 1830K 87M 18933 inet from any to any port = smtp queue	
15	P In L egress tcp K 31 5240 2 from <nospamd> to any port = smtp	

---

The `systat rules` view is especially useful because it offers a live view into the fully parsed and loaded rule set. For example, if your rule set behaves oddly, the rules view can point you in the right direction and show you the actual flow of packets.

The `systat` program also offers a view that presents the same data you would get via `pfctl -s` status on the command line. The following example shows part of the output of `systat pf`. The `systat pf` view offers more information than will fit on most screens, but you can page through the live display of the data.

---

2 users	Load 0.34 0.64 0.47 (1-16 of 51)	Fri Apr 1 14:04:04 2011
<hr/>		
TYPE NAME	VALUE	RATE NOTES
pf Status	Enabled	
pf Since	139:05:08	
pf Debug	err	
pf Hostid	0x82aea702	
nfe0 Bytes In	6217042900	IPv4
nfe0 Bytes In	0	IPv6
nfe0 Bytes Out	5993394114	IPv4
nfe0 Bytes Out	64	IPv6
nfe0 Packets In	12782504	IPv4, Passed
nfe0 Packets In	0	IPv6, Passed
nfe0 Packets In	11096	IPv4, Blocked
nfe0 Packets In	0	IPv6, Blocked
nfe0 Packets Out	12551463	IPv4, Passed
nfe0 Packets Out	1	IPv6, Passed
nfe0 Packets Out	167	IPv4, Blocked

---

The systat program offers quite a few other views, including network-related ones such as netstat, vmstat for virtual memory statistics, and iostat for input/output statistics by device. You can cycle through all systat views using the left and right cursor keys. (See `man systat` for full details.)

## ***Keeping an Eye on Things with pftop***

If your system does not have a have a systat version with the PF-related views, you can still keep an eye on what is passing in and out of your network in real time using Can Erkin Acar's `pftop`. This command shows a running snapshot of your traffic. `pftop` is not included in the base system, but it is available as a package, in ports on OpenBSD and FreeBSD as `sysutils/pftop`, and on NetBSD via `pkgsrc` as `sysutils/pftop`. Here is an example of its output:

---

pfTop: Up State 1-17/771, View: default, Order: none, Cache: 10000								14:05:42	
PR	DIR	SRC	DEST	STATE	AGE	EXP	PKTS	BYTES	
udp	Out	192.168.103.1:56729	192.168.103.9:12345	SINGLE:NO_TRAFFIC	8346m	22	373K	543M	
tcp	In	10.168.103.15:47185	213.187.179.198:22	ESTABLISHED:ESTABLISHED	62715	86395	3232	667K	
tcp	In	10.168.103.15:2796	213.187.179.198:22	ESTABLISHED:ESTABLISHED	62706	86369	4071	686K	
tcp	In	10.168.103.15:15599	129.240.64.10:6667	ESTABLISHED:ESTABLISHED	62336	86379	9318	854K	
tcp	Out	213.187.179.198:15599	129.240.64.10:6667	ESTABLISHED:ESTABLISHED	62336	86379	9318	854K	
tcp	In	10.168.103.15:8923	140.211.166.4:6667	ESTABLISHED:ESTABLISHED	62181	86380	15755	4821K	
tcp	Out	213.187.179.198:8923	140.211.166.4:6667	ESTABLISHED:ESTABLISHED	62181	86380	15755	4821K	
tcp	In	10.168.103.15:47047	217.17.33.10:6667	ESTABLISHED:ESTABLISHED	62146	86379	7132	559K	
tcp	Out	213.187.179.198:47047	217.17.33.10:6667	ESTABLISHED:ESTABLISHED	62146	86379	7132	559K	
tcp	In	10.168.103.15:30006	203.27.221.42:6667	ESTABLISHED:ESTABLISHED	62082	86380	6034	489K	
tcp	Out	213.187.179.198:30006	203.27.221.42:6667	ESTABLISHED:ESTABLISHED	62082	86380	6034	489K	
tcp	In	10.168.103.15:31709	209.250.145.51:6667	ESTABLISHED:ESTABLISHED	62082	86379	6685	617K	
tcp	Out	213.187.179.198:31709	209.250.145.51:6667	ESTABLISHED:ESTABLISHED	62082	86379	6685	617K	
tcp	In	192.168.103.254:53863	69.90.74.197:80	ESTABLISHED:ESTABLISHED	57056	29506	10	3282	
tcp	Out	213.187.179.198:53863	69.90.74.197:80	ESTABLISHED:ESTABLISHED	57056	29506	10	3282	
tcp	In	10.168.103.15:33241	192.168.103.84:22	ESTABLISHED:ESTABLISHED	47254	82340	7555	897K	
tcp	Out	10.168.103.15:33241	192.168.103.84:22	ESTABLISHED:ESTABLISHED	47254	82340	7555	897K	

---

You can use `pftop` to sort your connections by a number of different criteria, including by PF rule, volume, age, and source and destination addresses.

## ***Graphing Your Traffic with pfstat***

Once you have a system up and running and producing data, a graphical representation of traffic data is a useful way to view and analyze your data. One way to graph your PF data is with `pfstat`, a utility developed by Daniel Hartmeier to extract and present statistical data that is automatically generated by PF. The `pfstat` tool is available via the OpenBSD package system or as the port `net/pfstat`, via the FreeBSD ports system as `sysutils/pfstat`, and via NetBSD `pkgsrc` as `sysutils/pfstat`.

The `pfstat` program collects the data you specify in the configuration file and presents that data as JPG or PNG graphics files. The data source can be either PF running on the local system, via the `/dev/pf` device, or data collected from a remote computer running the companion `pfstated` daemon.

To set up *pfstat*, you simply decide which parts of your PF data you want to graph and how, and then write the configuration file and start cron jobs to collect the data and generate your graphs. The program comes with a well-annotated sample configuration file and a useful man page. The sample configuration is a useful starting point for writing your own configuration file. For example, the following *pfstat.conf* fragment is very close to one you will find in the sample configuration:<sup>1</sup>

---

```
collect 8 = global states inserts diff
collect 9 = global states removals diff
collect 10 = global states searches diff

image "/var/www/users/peter/bsdly.net/pfstat-states.jpg" {
    from 1 days to now
    width 980 height 300
    left
        graph 8 "inserts" "states/s" color 0 192 0 filled,
        graph 9 "removals" "states/s" color 0 0 255
    right
        graph 10 "searches" "states/s" color 255 0 0
}
```

---

The configuration here starts off with three collect statements, where each of the data series is assigned a unique numeric identifier. Here, we capture the number of insertions, removals, and searches in the state table. Next up is the image definition, which specifies the data that is to be graphed. The from line specifies the period to display (from 1 days to now means that only data collected during the last 24 hours is to be displayed). width and height specify the graph size measured in number of pixels in each direction. The graph statements specify how the data series are displayed and the graph legends. Collecting state insertions, removals, and searches once a minute, and then graphing the data collected over one day produces a graph roughly like the one in Figure 8-1.

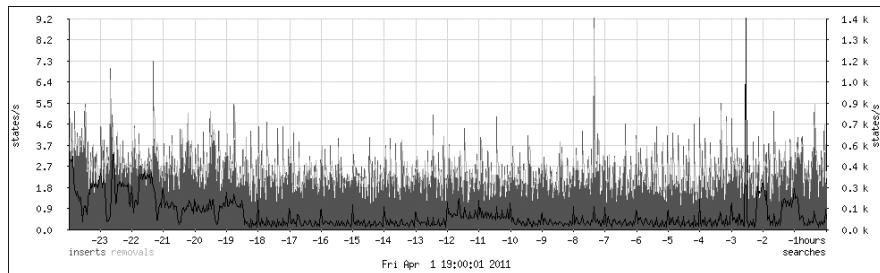


Figure 8-1: State table statistics, 24-hour time scale

---

1. The color values listed in the configuration example would give you a graph with red, blue, and green lines. For the print version of this book, we changed the colors to grayscale values: 0 192 0 became 105 105 105, 0 0 255 became 192 192 192, and 255 0 0 became 0 0 0.

The graph can be tweaked to provide a more detailed view of the same data. For example, to see the data for the last hour in a slightly higher resolution, change the period to from 1 hours to now and the dimensions to width 600 height 300. The result is something like the graph in Figure 8-2.

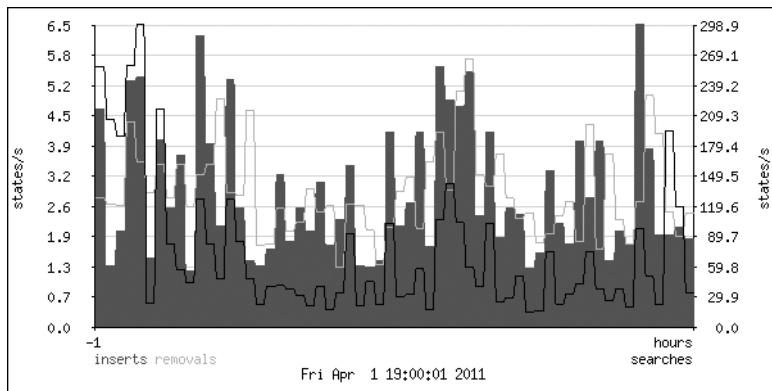


Figure 8-2: State table statistics, 1-hour time scale

The pfstat home page at <http://www.benzedrine.cx/pfstat.html> contains several examples, with demonstrations in the form of live graphs of the data from the *benzedrine.cx* domain's gateways. By reading the examples and tapping your own knowledge of your traffic, you should be able to create pfstat configurations that are well suited to your site's needs.

## Collecting NetFlow Data with pflow(4)

*NetFlow* is a network data collection and analysis method that has spawned many supporting tools for recording and analyzing data about TCP/IP connections. NetFlow originated at Cisco, and over time has become an essential feature in various network equipment as a tool for network management and analysis.

The NetFlow data model defines a network *flow* as a unidirectional sequence of packets with the same source and destination IP address and protocol. For example, a TCP connection will appear in NetFlow data as two flows: one in each direction.

PF data can be made available to NetFlow tools via the pflow(4) pseudo-interface that was introduced in OpenBSD 4.5 along with the pflow state option. Essentially, all the information you would expect to find in a NetFlow-style flow record is easily derived from the data PF keeps in the state table, and the pflow interface offers a straightforward way to export PF state table data in this processing-friendly and well-documented format. As with other logging, you enable NetFlow data collection in your PF rule set on a per-rule basis.

A complete NetFlow-based network monitoring system consists of several distinct parts. The NetFlow data originates at one or more *sensors* that generate data about network traffic. The sensors forward data about the flows to a

*collector*, which stores the data it receives. Finally, a *reporting* or *analysis* system lets you extract and process the data.<sup>2</sup>

## Setting Up the NetFlow Sensor

The NetFlow sensor requires two components: one or more configured pflow(4) devices and at least one pass rule in your rule set with the pflow state option enabled. The pflow interfaces are created with two required parameters: the flow source IP address and flow destination's IP address and port. Here is an example of the ifconfig command for the */etc/hostname.pflow0* file:

---

```
flowsrc 192.0.2.1 flowdst 192.0.2.105:3001
```

---

From the command line, use this command:

---

```
$ sudo ifconfig pflow0 create flowsrc 192.0.2.1 flowdst 192.0.2.105:3001
```

---

In both cases, this command sets up the host to send NetFlow data with a flow source address 192.0.2.1 to a collector that should listen for NetFlow data at 192.0.2.105, UDP port 3001.

**NOTE** *It is possible to set up several pflow devices with separate flow destinations. However, it is not currently possible to specify on a per-rule basis which pflow device should receive the generated data.*

After enabling the pflow device, specify in */etc/pf.conf* which pass rules should provide NetFlow data to the sensor. For example, if your main concern is to collect data on your clients' email traffic to IPv4 hosts, this rule would set up the necessary sensor:

---

```
pass out log inet proto tcp from <client> to port $email \
    label client-email keep state (pflow)
```

---

When pflow was first introduced to PF, the immediate reaction from early adopters was that more likely than not they would want to add the pflow option to most pass rules in their rule sets. This led PF developer Henning Brauer to introduce another useful PF feature, the ability to set *state defaults* that apply to all rules unless otherwise specified. For example, if you add the following line at the start of your rule set, all pass rules in the configuration will generate NetFlow data to be exported via the pflow device.

---

```
set state-defaults pflow
```

---

With at least one pflow device configured and at least one rule in your *pf.conf* that generates data for export via the pflow device, you're almost finished setting up the sensor. However, you may still need to add a rule that allows the UDP data to flow from the IP address you specified as the flow data

---

2. For a more in-depth treatment of network analysis with NetFlow-based tools, see *Network Flow Analysis* by Michael W. Lucas (No Starch Press, 2010).

source to the collector's IP address and target port at the flow destination. Once you've completed this last step, you should be ready to turn your attention to collecting the data for further processing.

## NetFlow Data Collecting, Reporting, and Analysis

If your site has a NetFlow-based collection and analysis infrastructure in place, you may already have added the necessary configuration to feed the PF-originated data into the data collection and analysis system. If you have not yet set up a flow-analysis environment, there are a number of options available.

The OpenBSD packages system offers three NetFlow collector and analysis packages: `flow-tools`, `flowd`, and `nfdump`.<sup>3</sup> All three systems have a dedicated and competent developer and user community, and various addons, including graphical web interfaces. `flow-tools` is the main component in many sites' flow-analysis setups. The `nfdump` fans point to the `nfsen` analysis package that integrates the `nfdump` tools in a powerful and flexible web-based analysis front end.

### CHOOSING A COLLECTOR

The choice of collector is somewhat tied to the choice of analysis package. Perhaps because the collectors tend to store flow data in their own unique formats, most reporting and analysis back ends are developed with a distinctive bias for one or the other collector.

Regardless of your choice of NetFlow collector, the familiar logging caveats apply: Detailed traffic log information will require storage. In the case of NetFlow, each flow will generate a record of fairly fixed size, and anecdotal evidence indicates that even modest collection profiles on busy sites can generate gigabytes of NetFlow data per day. The amount of storage you'll need is directly proportional to the number of connections and how long you keep the original NetFlow data. Finally, recording and storing traffic logs with this level of detail is likely to have legal implications.

Collectors generally offer filtering features that let you discard data about specific hosts or networks, or even discard some parts of the NetFlow records themselves, either globally or for data about specific hosts or networks.

To illustrate some basic NetFlow collection and how to extract a subset of the collected data for further analysis, we'll use `flowd`, developed by long-time OpenBSD developer Damien Miller and available via the package systems (on OpenBSD as `net/flowd` and on FreeBSD as `net-mgmt/flowd`).

I've chosen to use `flowd` here mainly because it was developed to be small, simple, and secure. As you will see, `flowd` still manages to be quite useful and flexible. Flow data operations with other tools will differ in some details, but the underlying principles remain the same.

<sup>3</sup> The actively maintained project home pages for `flow-tools` and `nfdump` are <http://code.google.com/p/flow-tools/> and <http://nfdump.sourceforge.net/>. (The older versions should still be available from <http://www.splintered.net/sw/flow-tools/>.) The `nfsen` web front end has a project page at <http://nfsen.sourceforge.net/>. For the latest information about `flowd`, visit <http://www.mindrot.org/flowd.html>.

When compared to other NetFlow collector suites, `flowd` is very compact, with only two executable programs (the collector daemon `flowd` and the flow-filtering and presentation program `flowd-reader`), as well as the supporting library and controlling configuration file. The documentation is adequate if a bit terse, and the sample `/etc/flowd.conf` file contains a generous amount of comments. Based on the man pages and the comments in the sample configuration file, it should not take you long to create a useful collector configuration.

After stripping out any comment lines (using `grep -v \# /etc/flowd.conf` or similar), a very basic `flowd` configuration could look like this:

---

```
logfile "/var/log/flowd"
listen on 192.0.2.105:3001
flow source 192.0.2.1
store ALL
```

---

While this configuration barely contains more information than the `pflow` interface's configuration in the earlier description of setting up the sensor, it does include two important items:

- The `logfile` line tells us where the collected data is to be stored (and reveals that `flowd` tends to store all data in a single file).
- The final line tells us that `flowd` will store all fields in the data it receives from the designated flow source.

With this configuration in place, start up the `flowd` daemon, and almost immediately you should see the `/var/log/flowd` file grow as network traffic passes through your gateway and flow records are collected. After a while, you should be able to look at the data using `flowd`'s companion program `flowd-reader`. For example, with all fields stored, the data for one name lookup from a host on the NATed local network looks like this in `flowd-reader`'s default view:

---

```
$ sudo flowd-reader /var/log/flowd
FLOW recv_time 2011-04-01T21:15:53.607179 proto 17 tcpflags 00 tos 00 agent
[213.187.179.198] src [192.168.103.254]:55108 dst [192.168.103.1]:53 packets 1
octets 62
FLOW recv_time 2011-04-01T21:15:53.607179 proto 17 tcpflags 00 tos 00 agent
[213.187.179.198] src [192.168.103.1]:53 dst [192.168.103.254]:55108 packets 1
octets 129
```

---

Notice that the lookup generates two flows: one in each direction.

The first flow is identified mainly by the time it was received, followed by the protocol used (protocol 17 is UDP, as `/etc/protocols` will tell you). The connection had both TCP and TOS flags unset, and the collector received the data from our gateway at 192.0.2.1. The flow's source address was 192.168.103.254, source port 55108, and the destination address was 192.168.103.1 and port 53, conventionally the DNS port. The flow consisted of one packet, with a payload of 62 octets. The return flow was received by the collector at the same time, and we see that this flow has the source and destination reversed, with a

slightly larger payload of 129 octets. `flowd-reader`'s output format lends itself to parsing by regular expressions for postprocessing in reporting tools or plotting software.

You might think that this data is all anyone would ever want to know about any particular set of network flows, but it is possible to extract even more detailed information. For example, using the `flowd-reader -v` option for verbose output, you might see something like this:

---

```
FLOW recv_time 2011-04-01T21:15:53.607179 proto 17 tcpflags 00 tos 00 agent [213.187.179.198] src [192.168.103.254]:55108 dst [192.168.103.1]:53 gateway [0.0.0.0] packets 1 octets 62 in_if 0 out_if 0 sys_uptime_ms 1w5d19m59s.000 time_sec 2011-04-01T21:15:53 time_nanosec 103798508 netflow ver 5 flow_start 1w5d19m24s.000 flow_finish 1w5d19m29s.000 src_AS 0 src_masklen 0 dst_AS 0 dst_masklen 0 engine_type 10752 engine_id 10752 seq 5184351 source 0 crc32 759adcbdb
FLOW recv_time 2011-04-01T21:15:53.607179 proto 17 tcpflags 00 tos 00 agent [213.187.179.198] src [192.168.103.1]:53 dst [192.168.103.254]:55108 gateway [0.0.0.0] packets 1 octets 129 in_if 0 out_if 0 sys_uptime_ms 1w5d19m59s.000 time_sec 2011-04-01T21:15:53 time_nanosec 103798508 netflow ver 5 flow_start 1w5d19m24s.000 flow_finish 1w5d19m29s.000 src_AS 0 src_masklen 0 dst_AS 0 dst_masklen 0 engine_type 10752 engine_id 10752 seq 5184351 source 0 crc32 f43ccb22
```

---

The `gateway` field indicates that the sensor itself served as the gateway for this connection. You see a list of the interfaces involved (the `in_if` and `out_if` values), the sensor's system uptime (`sys_uptime_ms`), and a host of other parameters such as AS numbers (`src_AS` and `dst_AS`) that may be useful for statistics or filtering purposes in various contexts. Once again, the output is ideally suited to filtering via regular expressions.

You do not need to rely on external software for the initial filtering on the data you collect from your pflow sensor. `flowd` itself offers a range of filtering features that make it possible to store only the data you need. One approach is to put the filtering expressions in the `flowd.conf`, as in the following example (with the comments stripped to save space):

---

```
logfile "/var/log/flowd.compact"
listen on 192.0.2.105:3001
flow source 192.0.2.1
store SRC_ADDR
store DST_ADDR
store SRCDST_PORT
store PACKETS
store OCTETS
internalnet = "192.168.103.0/24"
unwired = "10.168.103.0/24"
discard src $internalnet
discard dst $internalnet
discard src $unwired
discard dst $unwired
```

---

You can choose to store only certain fields in the flow records. For example, in configurations where there is only one collector or agent, the agent field serves no useful purpose and does not need to be stored. In this configuration, we choose to store only the source and destination address and port, the number of packets, and the number of octets.

You can limit the data you store even further. The macros `internalnet` and `unwired` expand to two NATed local networks, and the four `discard` lines following the macro definitions mean that `flowd` discards any data it receives about flows with either source or destination addresses in either of those local networks. The result is a more compact set of data, tailored to your specific needs, and you see only routable addresses and the address of the sensor gateway's external interface:

---

```
$ sudo flowd-reader /var/log/flowd.compact | head
FLOW src [193.213.112.71]:38468 dst [192.0.2.1]:53 packets 1 octets 79
FLOW src [192.0.2.1]:53 dst [193.213.112.71]:38468 packets 1 octets 126
FLOW src [200.91.75.5]:33773 dst [192.0.2.1]:53 packets 1 octets 66
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:33773 packets 1 octets 245
FLOW src [200.91.75.5]:3310 dst [192.0.2.1]:53 packets 1 octets 75
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:3310 packets 1 octets 199
FLOW src [200.91.75.5]:2874 dst [192.0.2.1]:53 packets 1 octets 75
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:2874 packets 1 octets 122
FLOW src [192.0.2.1]:15393 dst [158.37.91.134]:123 packets 1 octets 76
FLOW src [158.37.91.134]:123 dst [192.0.2.1]:15393 packets 1 octets 76
```

---

Even with the `-v` option, `flowd-reader`'s display reveals only what you explicitly specify in the filtering configuration:

---

```
$ sudo flowd-reader -v /var/log/flowd.compact | head
LOGFILE /var/log/flowd.compact
FLOW src [193.213.112.71]:38468 dst [192.0.2.1]:53 packets 1 octets 79
FLOW src [192.0.2.1]:53 dst [193.213.112.71]:38468 packets 1 octets 126
FLOW src [200.91.75.5]:33773 dst [192.0.2.1]:53 packets 1 octets 66
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:33773 packets 1 octets 245
FLOW src [200.91.75.5]:3310 dst [192.0.2.1]:53 packets 1 octets 75
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:3310 packets 1 octets 199
FLOW src [200.91.75.5]:2874 dst [192.0.2.1]:53 packets 1 octets 75
FLOW src [192.0.2.1]:53 dst [200.91.75.5]:2874 packets 1 octets 122
FLOW src [192.0.2.1]:15393 dst [158.37.91.134]:123 packets 1 octets 76
```

---

Fortunately, `flowd` does not force you to make all your filtering decisions when your collector receives the flow data from the sensor. Using the `-f` flag, you can specify a separate file with filtering statements to extract specific data from a larger set of collected flow data. For example, to see HTTP traffic to your web server, you could write a filter that stores only flows with your web server's address and TCP port 80 as the destination, or your web server and TCP port 80 as the source:

---

```
webserver = 192.0.2.227
discard all
accept dst $webserver port 80 proto tcp
```

---

```
accept src $webserver port 80 proto tcp
store RECV_TIME
store SRC_ADDR
store DST_ADDR
store PACKETS
store OCTETS
```

---

Assuming you stored the filter in `towebserver.flowdfilter`, you could then extract traffic matching your filtering criteria from `/var/log/flowd` like this:

```
$ sudo flowd-reader -v -f towebserver.flowdfilter /var/log/flowd | tail
FLOW recv_time 2011-04-01T21:13:15.505524 src [89.250.115.174] dst
[192.0.2.227] packets 6 octets 414
FLOW recv_time 2011-04-01T21:13:15.505524 src [192.0.2.227] dst
[89.250.115.174] packets 4 octets 725
FLOW recv_time 2011-04-01T21:13:49.605833 src [216.99.96.53] dst [192.0.2.227]
packets 141 octets 7481
FLOW recv_time 2011-04-01T21:13:49.605833 src [192.0.2.227] dst [216.99.96.53]
packets 212 octets 308264
FLOW recv_time 2011-04-01T21:14:04.606002 src [91.121.94.14] dst [192.0.2.227]
packets 125 octets 6634
FLOW recv_time 2011-04-01T21:14:04.606002 src [192.0.2.227] dst [91.121.94.14]
packets 213 octets 308316
FLOW recv_time 2011-04-01T21:14:38.606384 src [207.46.199.44] dst
[192.0.2.227] packets 10 octets 642
FLOW recv_time 2011-04-01T21:14:38.606384 src [192.0.2.227] dst
[207.46.199.44] packets 13 octets 16438
FLOW recv_time 2011-04-01T21:15:14.606768 src [213.187.176.94] dst
[192.0.2.227] packets 141 octets 7469
FLOW recv_time 2011-04-01T21:15:14.606768 src [192.0.2.227] dst
[213.187.176.94] packets 213 octets 308278
```

---

In addition to the filtering options demonstrated here, the `flowd` filtering functions take a number of other options. Some of those options will be quite familiar from other filtering contexts such as PF, including a range of network-oriented parameters; others are more oriented to extracting data on flows originating at specific dates or time periods and other storage-oriented parameters. The full story, as always, is found in `man flowd.conf`.

Once you have extracted the data you need, you have several tools available for processing and presenting your data.

## **Collecting NetFlow Data with `pfflowd`**

For systems that do not support NetFlow data export via `pflogd`, NetFlow support is available via the `pfflowd` package. The NetFlow data model defines a network flow as a unidirectional sequence of packets with the same source and destination IP address, and protocol. This maps very well to PF state information, and `pfflowd` is intended to record state changes from the local system's `pfsync` device. Once enabled, `pfflowd` acts as a NetFlow sensor that converts `pfsync` data to NetFlow format for transmission to a NetFlow collector on the network.

The `pfflowd` tool was written and is maintained by Damien Miller, and is available from <http://www.mindrot.org/projects/pfflowd/>, as well as through the package systems on OpenBSD and FreeBSD as `net/pfflowd`. The lack of `pfsync` support on NetBSD means that `pfflowd` is not available on that platform as of this writing.

### **SNMP Tools and PF-Related SNMP MIBs**

Simple Network Management Protocol (SNMP) was designed to let network administrators collect and monitor key data about how their systems run, and change configurations on multiple network nodes from a centralized system.<sup>4</sup> The SNMP protocol comes with a well-defined interface and a method for extending the Management Information Base (MIB), which defines the managed devices and objects.

Both proprietary and open source network management and monitoring systems generally have SNMP support in one form or the other, and in some products, it's a core feature. On the BSDs, SNMP support has generally come in the form of the `net-snmp` package, which provides the tools you need to retrieve SNMP data and to collect data for retrieval by management systems. The package is available on OpenBSD as `net/net-snmp`, on FreeBSD as `net-mgmt/net-snmp`, and on NetBSD as `net/net-snmp`. OpenBSD's `snmpd` (written mainly by Reyk Floeter) debuted as part of the base system in OpenBSD 4.3, and implements all required SNMP functionality. (See `man snmpd` and `man snmpd.conf` for details.)

There are MIBs to make PF data available to SNMP monitoring. Joel Knight maintains the MIBs for retrieving data on PF, CARP, and OpenBSD kernel sensors, and offers them for download from <http://www.packetmischief.ca/openbsd/snmp/>. The site also offers patches to the `net-snmp` package to integrate the OpenBSD MIBs.

After installing the package and the extension, your SNMP-capable monitoring systems will be able to watch PF data in any detail you desire. (FreeBSD's `bsnmpd` includes a PF module. See the `bsnmpd` man page for details.)

## **Log Data as the Basis for Effective Debugging**

In this chapter, we walked through the basics of collecting, displaying, and interpreting data about a running system with PF enabled. Knowing how to find and use information about how your system behaves is useful for several purposes.

Keeping track of the status of a running system is useful in itself, but the ability to read and interpret log data is even more essential when testing your setup. Another prime use for log data is for tracking the effect of changes you make in the configuration, such as when tuning your system to give optimal performance. In the next chapter, we'll focus on checking your configuration and tuning it for optimal performance, based on log data and other observations.

---

4. The protocol debuted with RFC 1067 in August 1988 and is now in its third major version as defined in RFCs 3411 through 3418.

# 9

## GETTING YOUR SETUP JUST RIGHT



By now, you have spent significant time designing your network and implementing that design in your PF configuration. Getting your setup just right—removing any remaining setup bugs and inefficiencies—can be quite challenging at times.

This chapter describes some options and methods that will help you get the setup you need. First, we will take a look at global options and some settings that can have a profound influence on how your configuration behaves.

### Things You Can Tweak and What You Probably Should Leave Alone

Network configurations are inherently very tweakable. While browsing the `pf.conf` man page or other reference documentation, it is easy to be overwhelmed by the number of options and settings that you could conceivably adjust in order to get that perfectly optimized setup.

Keep in mind that for PF in general, *the defaults are sane* for most setups. Some settings and variables lend themselves to tuning; others should come with a big warning that they should be adjusted only in highly unusual circumstances, if at all.

Here, we'll look at some of the global settings that you should know about, although you won't need to change them in most circumstances.

These options are written as `set option setting` and go after any macro definitions in your `pf.conf` file, but before translation or filtering rules.

**NOTE** *If you read the pf.conf man page, you will discover that a few other options are available. However, most of those are not relevant in a network testing and performance tuning context.*

## **Block Policy**

The `block-policy` option determines which feedback, if any, PF will give to hosts that try to create connections that are subsequently blocked. The option has two possible values:

- `drop` drops blocked packets with no feedback.
- `return` returns with status codes such as `Connection refused` or similar.

The correct strategy for block policies has been the subject of considerable discussion over the years. The default setting for `block-policy` is `drop`, which means that the packet is silently dropped without any feedback. However, silently dropping packets makes it likely that the sender will resend the unacknowledged packets, rather than drop the connection. Thus, the effort is kept up until the relevant timeout counter expires. Unless you can think of a good reason to set the block policy to something else, set it to `return`:

---

```
set block-policy return
```

---

This means that the sender's networking stack will receive an unambiguous signal indicating that the connection was refused.

This setting specifies the *global* default for your block policy. If necessary, you can still vary the blocking type for specific rules.

For example, you could change the brute-force protection rule set from Chapter 6 to have `block-policy` set to `return`, but use `block drop quick from <bruteforce>` to make the brute forceers waste time if they stick around once they have been added to the `<bruteforce>` table. You could also specify `drop` for traffic from nonroutable addresses coming in on your Internet-facing interface.

## **Skip Interfaces**

The `skip` option lets you exclude specific interfaces from all PF processing. The net effect is like a pass-all rule for the interface, but actually disables all PF processing on the interface. One common example is to disable filtering on

the loopback interface group, where filtering in most configurations adds little in terms of security or convenience, as follows:

---

```
set skip on lo
```

---

In fact, filtering on the loopback interface is almost never useful, and can lead to odd results with a number of common programs and services. The default is that `skip` is unset, which means that all configured interfaces can take part in PF processing. In addition to making your rule set slightly simpler, setting `skip` on interfaces where you do not want to perform filtering results in a slight performance gain.

## **State Policy**

The `state-policy` option specifies how PF matches packets to the state table. There are two possible values:

- With the default `floating` state policy, traffic can match state on all interfaces, not just the one where the state was created.
- With an `if-bound` policy, traffic will match only on the interface where the state is created; traffic on other interfaces will not match the existing state.

Like `block-policy`, this option specifies the global state-matching policy.

You can override state policy on a per-rule basis if needed. For example, in a rule set with the default `floating` state policy, you could have a rule like this:

---

```
pass out on egress inet proto tcp to any port $allowed modulate state (if-bound)
```

---

With this rule, any return traffic trying to pass back in would need to pass on the same interface where the state was created in order to match the state-table entry.

The situations in which `state-policy if-bound` is useful are rare enough that you should leave this setting at the default.

## **State Defaults**

The `state-defaults` option was introduced in OpenBSD 4.5 to enable setting specific state options as the default options to all rules in the rule set unless specifically overridden by other options in individual rules.

Here is a common example:

---

```
set state-defaults pflow
```

---

This sets up all `pass` rules in the configuration to generate NetFlow data to be exported via a `pflow` device.

In some contexts, it makes sense to apply state-tracking options, such as connection limits, as a global state default for the entire rule set. Here is an example:

---

```
set state-defaults max 1500, max-src-conn 100, source-track rule
```

---

This sets the default maximum number of state entries per rule to 1,500, with a maximum of 100 simultaneous connections from any one host, with separate limits for each rule in the loaded rule set.

Any option that is valid inside parentheses for `keep state` in an individual rule can also be included in a `set state-defaults` statement. Setting state defaults in this way is useful if there are state options that are not already system defaults that you want to apply to all rules in your configuration.

## Timeouts

The `timeout` option sets the timeouts and related options for various interactions with the state-table entries. The majority of the available parameters are protocol-specific values stored in seconds and prefixed `tcp.`, `udp.`, `icmp.`, and `other..`. However, `adaptive.start` and `adaptive.end` denote the number of state-table entries.

The following `timeout` options affect state-table memory use and to some extent lookup speed:

- The `adaptive.start` and `adaptive.end` values set the limits for scaling down timeout values once the number of state entries reach the `adaptive.start` value. When the number of states reaches `adaptive.end`, all timeouts are set to 0, essentially expiring all states immediately. The defaults are 6,000 and 12,000 (calculated as 60 percent and 120 percent of the state limit), respectively. These settings are intimately related to the memory pool limit parameters you set via the `limit` option.
- The `interval` value denotes the number of seconds between purges of expired states and fragments. The default is 10 seconds.
- The `frag` value denotes the number of seconds a fragment will be kept in an unassembled state before it is discarded. The default is 30 seconds.
- When set, `src.track` denotes the number of seconds source-tracking data will be kept after the last state has expired. The default is 0 seconds.

You can inspect the current settings for all `timeout` parameters with `pfctl -s timeouts`. For example, the following display shows a system running with default values.

---

```
$ sudo pfctl -s timeouts
tcp.first          120s
tcp.opening         30s
tcp.established    86400s
tcp.closing         900s
tcp.finwait        45s
tcp.closed          90s
```

---

---

tcp.tsdiff	30s
udp.first	60s
udp.single	30s
udp.multiple	60s
icmp.first	20s
icmp.error	10s
other.first	60s
other.single	30s
other.multiple	60s
frag	30s
interval	10s
adaptive.start	6000 states
adaptive.end	12000 states
src.track	0s

---

These options can be used to tweak your setup for performance. However, changing the protocol-specific settings from the default values creates a significant risk that valid but idle connections might be dropped prematurely or blocked outright.

## Limits

The `limit` option sets the size of the memory pools PF uses for state tables and address tables. These are hard limits, so you may need to increase or tune the values for various reasons. If your network is a busy one with larger numbers than the default values allow for, or if your setup requires large address tables or a large number of tables, then this section will be very relevant for you.

Keep in mind that the total amount of memory available through memory pools is taken from the *kernel memory space*, and the total available is a function of total available kernel memory. The kernel allocates a fixed amount of memory for its own use at system startup. However, since kernel memory is never swapped, the amount of memory allocated to the kernel can never equal or exceed all physical memory in the system. (If that happened, there would be no space for user mode programs to run.)

The amount of available pool memory depends on which hardware platform you use, as well as on a number of hard-to-predict variables specific to the local system. On the i386 architecture, the maximum kernel memory is in the 768MB to 1GB range, depending on a number of factors, including the number and kind of hardware devices in the system. The amount actually available for allocation to memory pools comes out of this total, again depending on a number of system-specific variables.

To inspect the current limit settings, use `pfctl -sm`. Typical output looks like this:

---

```
$ sudo pfctl -sm
states      hard limit    10000
src-nodes   hard limit    10000
frags       hard limit    5000
tables      hard limit    1000
table-entries hard limit  200000
```

---

To change these values, edit *pf.conf* to include one or more lines with new limit values. For example, you could use the following lines to raise the hard limit for number of states to 25,000 and table entries to 300,000:

---

```
set limit states 25000
set limit table-entries 300000
```

---

You can also set several limit parameters at the same time in a single line by enclosing them in brackets, like this:

---

```
set limit { states 25000, src-nodes 25000, table-entries 300000 }
```

---

In the end, you almost certainly should not change the limits at all. If you do, however, it is important to watch your system logs for any indication that your changed limits do not have undesirable side effects or do not fit in available memory. Setting the debug level to a higher value is potentially quite useful for watching the effects of tuning limit parameters.

## Debug

The debug option determines what, if any, error information PF will generate at the *kern.debug* log level. The default value is err, which means that only serious errors will be logged. Since OpenBSD 4.7, the log levels here correspond to the ordinary syslog levels, which range from emerg (panics are logged), alert (correctable but very serious errors logged), crit (critical conditions logged), err (errors are logged), warning (warnings are logged), notice (unusual conditions are logged), info (informational messages are logged), to debug (full debugging information, likely only useful to developers).

In pre-OpenBSD 4.7 versions, PF used its own log-level system, with a default of urgent (equivalent to err in the new system). The other possible settings were none (no messages), misc (reporting slightly more than urgent), and loud (producing status messages for most operations). The *pfctl* parser still accepts the older-style debug levels for compatibility.

After running one of my gateways at the debug level for a little while, this is what a typical chunk of the */var/log/messages* file looked like:

---

```
$ tail -f /var/log/messages
Oct  4 11:41:11 skapet /bsd: pf_map_addr: selected address 194.54.107.19
Oct  4 11:41:15 skapet /bsd: pf: loose state match: TCP 194.54.107.19:25
194.54.107.19:25 158.36.191.135:62458 [lo=3178647045 high=3178664421 win=33304
modulator=0 wscale=1] [lo=3111401744 high=3111468309 win=17376 modulator=0
wscale=0] 9:9 R seq=3178647045 (3178647044) ack=3111401744 len=0 ackskew=0
pkts=9:12
Oct  4 11:41:15 skapet /bsd: pf: loose state match: TCP 194.54.107.19:25
194.54.107.19:25 158.36.191.135:62458 [lo=3178647045 high=3178664421 win=33304
modulator=0 wscale=1] [lo=3111401744 high=3111468309 win=17376 modulator=0
wscale=0] 10:10 R seq=3178647045 (3178647044) ack=3111401744 len=0 ackskew=0
pkts=10:12
Oct  4 11:42:24 skapet /bsd: pf_map_addr: selected address 194.54.107.19
```

---

As you can see, the debug level gives you a level of detail where PF repeatedly reports the IP address for the interface it is currently handling. In between the selected address messages, PF warns twice for the same packet that the sequence number is at the very edge of the expected range. This level of detail seems almost breathtaking at first glance, but in some circumstances, studying this kind of output is the best way to diagnose a problem and later check to see if your solution actually helped.

**NOTE** *This option can be set from the command line with pfctl -x, followed by the debug level you want. The command pfctl -x debug gives you maximum debugging information; pfctl -x none turns off debug messages entirely.*

Keep in mind that some debug settings can produce large amounts of log data, and in extreme cases, could impact performance all the way to self-denial-of-service level.

## **Rule Set Optimization**

The ruleset-optimization option enables or sets the mode for the rule set optimizer. The default setting for ruleset-optimization in OpenBSD 4.1 and equivalents is none, which means that no rule set optimization is performed at load time. From OpenBSD 4.2 onward, the default is basic, which means that when the rule set loads, the optimizer does the following things:

- Removes duplicate rules
- Removes rules that are subsets of other rules
- Merges rules into tables if appropriate (typical rule-to-table optimizations are rules that pass, redirect, or block based on identical criteria except source and/or target addresses)
- Changes the order of rules to improve performance

For example, say you have the macro `tcp_services = { ssh, www, https }` combined with the rule `pass proto tcp from any to self port $tcp_services`. Elsewhere in your rule set, you have a different rule that says `pass proto tcp from any to self port ssh`. The second rule is clearly a subset of the first, and they can be merged into one. Another common combination is having a `pass` rule like `pass proto tcp from any to int_if:network port $tcp_services` with otherwise identical `pass` rules where the target addresses are all in the `int_if:network` range.

With ruleset-optimization set to profile, the optimizer analyzes the loaded rule set relative to the actual network traffic in order to determine the optimal order of quick rules.

You can also set the value of the optimization option from the command line with `pfctl`:

---

```
$ sudo pfctl -o basic
```

---

This example enables the rule set optimization in basic mode.

Since the optimization may remove or reorder rules, the meaning of some statistics—mainly the number of evaluations per rule—may change in ways that may be hard to predict. In most cases, however, the effect is negligible.

### ***Optimization***

The `optimization` option specifies profiles for state-timeout handling. The possible values are `normal`, `high-latency`, `satellite`, `aggressive`, and `conservative`. The recommendation is to keep the default `normal` setting unless you have very specific needs.

The values `high-latency` and `satellite` are synonyms, where states expire more slowly in order to compensate for potential high latency.

The `aggressive` setting expires states early in order to save memory. This could, in principle, increase the risk of dropping idle-but-valid connections if your system is already close to its load and traffic limits, but anecdotal evidence indicates that the `aggressive` optimization setting rarely, if ever, interferes with valid traffic.

The `conservative` setting goes to great lengths to preserve states and idle connections, at the cost of some additional memory use.

### ***Fragment Reassembly***

The fragment reassembly options tied to `scrub` were significantly reworked in OpenBSD 4.6, introducing the new `set reassemble` option to turn reassembly of fragmented packets on or off. If `reassemble` is set to `off`, fragmented packets are simply dropped. The default is `set reassemble on`, which means fragments are reassembled and reassembled packets where the `do not fragment` bit was set on individual fragments will have the bit cleared.

## **Cleaning Up Your Traffic**

The next two features we'll discuss, `scrub` and `antispoof`, share a common theme: They provide automated protection against potentially dangerous clutter in your network traffic. Together, they are commonly referred to as tools for “network hygiene,” because they sanitize your networking considerably.

### ***Packet Normalization with scrub***

In PF versions up to OpenBSD 4.5 inclusive, the `scrub` keyword enables network traffic normalization. With `scrub`, fragmented packets are reassembled, and invalid fragments such as overlapping fragments are discarded, so the resulting packet is complete and unambiguous.

Enabling scrub provides a measure of protection against certain kinds of attacks based on incorrect handling of packet fragments.<sup>1</sup> A number of supplementing options are available, but the simplest form is suitable for most configurations.

---

`scrub in`

---

In order for some services to work with scrub, specific options must be set. One commonly cited example is some NFS implementations, which will not work with scrub at all unless you use the `no-df` parameter to clear the do not fragment bit on any packets that have the bit set. Some combinations of services, operating systems, and network configurations may require some of the more exotic scrub options.

In OpenBSD 4.6, scrub was demoted from stand-alone rule material to become an action you could attach to `pass` or `match` rules (the introduction of `match` rules being one of the main new PF features in OpenBSD 4.6). One other important development in the same rewrite of the scrub code was that the numerous packet-reassembly options were eliminated in favor of the new `reassemble` option, which simply turns reassembly on or off.

With the new scrub syntax, you need to supply at least one option in parentheses. The following works quite well for several networks in my care:

---

`match in all scrub (no-df max-mss 1440)`

---

This clears the do not fragment bit and sets the maximum segment size to 1440 bytes.

Other variations are possible, and even though the list of scrub options shrank somewhat for the OpenBSD 4.6 version, you should be able to cater to various specific needs by consulting the man pages and doing some experimentation.

If you find yourself needing to debug a scrub-related problem, study the `pf.conf` man page and consult the gurus on the relevant mailing lists.

## ***Protecting Against Spoofing with antispoof***

There are some very useful and common packet-handling actions that could be written as PF rules, but not without becoming long, complicated, and error-prone rule set boilerplate. Thus, `antispoof` was implemented for a common special case of filtering and blocking. This mechanism protects against activity from spoofed or forged IP addresses, mainly by blocking packets that appear on interfaces traveling in directions that are not logically possible.

---

1. Some notable attack techniques, including several historical denial-of-service setups, have exploited bugs in fragment handling that could lead to out-of-memory conditions or other resource exhaustion. One such exploit, which was aimed at Cisco's PIX firewall series, is described in the advisory at [http://www.cisco.com/en/US/products/products\\_security\\_advisory09186a008011e78d.shtml](http://www.cisco.com/en/US/products/products_security_advisory09186a008011e78d.shtml).

With `antispoof`, you can specify that you want to weed out spoofed traffic coming in from the rest of the world and any spoofed packets that (however unlikely) were to originate in your own network. Figure 9-1 illustrates the concept.

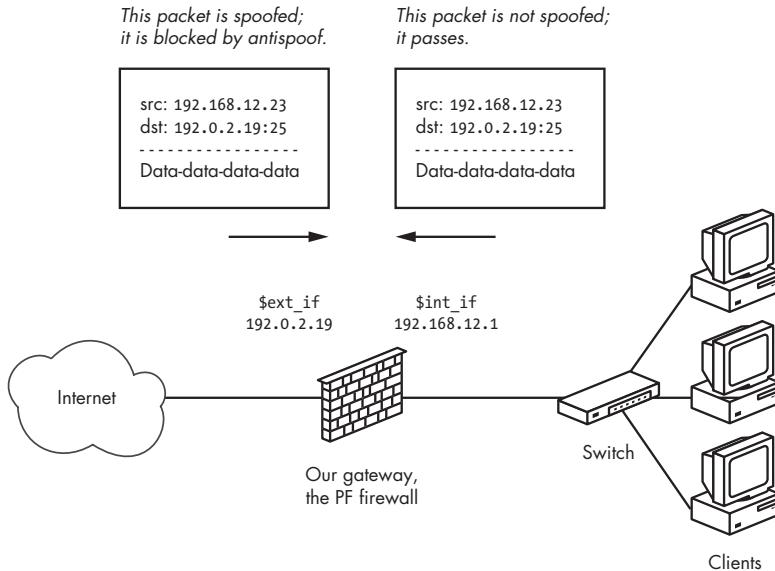


Figure 9-1: *antispooft drops packets that come in from the wrong network.*

To establish the kind of protection depicted in the diagram, specify `antispooft` for both interfaces in the illustrated network with these two lines:

---

```
antispooft for $ext_if
antispooft for $int_if
```

---

These lines expand to complex rules. The first one blocks incoming traffic when the source address appears to be part of the network directly connected to the antispoofed interface but arrives on a different interface. The second rule performs the same functions for the internal interface, blocking any traffic with apparently local network addresses that arrive on other interfaces than `$int_if`. However, keep in mind that `antispooft` is not designed to detect address spoofing remotely for networks that are not directly connected to the machine running PF.

## Testing Your Setup

Now it's time to dust off the precise specification that describes how your setup *should* work.

The physical layout of our sample network is centered around a *gateway* connected to the Internet via `$ext_if`. Attached to the gateway via `$int_if` is a *local network* with workstations and possibly one or more servers for local use.

Finally, we have a *DMZ* connected to `$dmz_if`, populated with servers offering services to the local network and the Internet. Figure 9-2 shows the logical layout of the network.

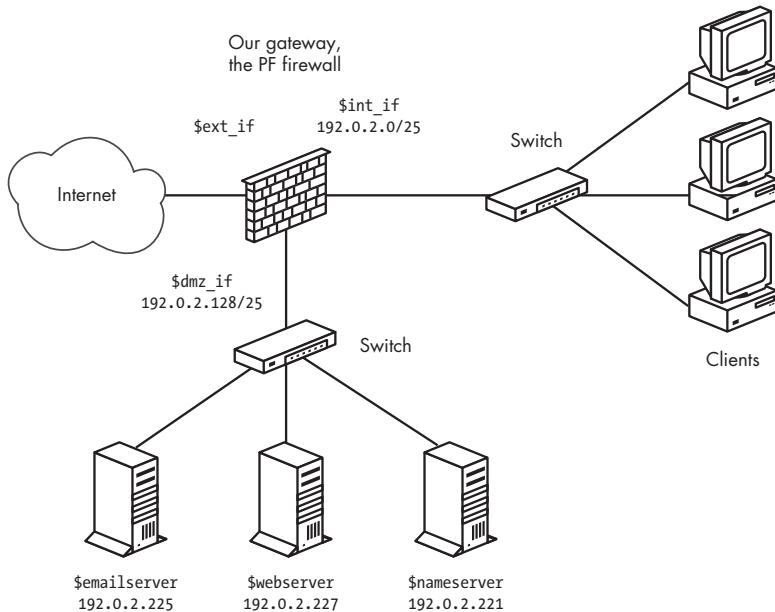


Figure 9-2: Network with servers in a DMZ

The corresponding rule set specification looks something like this:

- Machines outside our network should have access to the services offered by our servers in the DMZ, and no access to the local network.
- The machines in our local network, attached to `$int_if`, should have access to the services offered by the servers in the DMZ and access to a defined list of services outside our network.
- The machines in the DMZ should have access to some network services in the outside world.

The task at hand is to make sure the rule set we have in place actually implements the specification. We need to test the setup. A useful test would be to try the sequence in Table 9-1.

Your configuration may call for other tests or could differ in some particulars. Your real-life test scenario should specify how packets and connections should be logged. The main point is that you should decide what the expected and desired result for each of your test cases should be before you start testing.

In general, you should test using the applications you expect the typical user to have, such as web browsers or mail clients on various operating systems. The connections should simply succeed or fail, according to specifications. If one or more of your basic tests gives an unexpected result, move on to debugging your rule set.

**Table 9-1:** Sample Rule Set Test Case Sequence

Test Action	Expected Result
Try a connection from the local network to each allowed port on the servers in the DMZ.	The connection should pass.
Try a connection from the local network to each allowed port on servers outside our network.	The connection should pass.
Try a connection on any port from the DMZ to the local network.	The connection should be blocked.
Try a connection from the DMZ to each allowed port on servers outside our network.	The connection should pass.
Try a connection from outside our network to \$webserver in the DMZ on each port in \$webports.	The connection should pass.
Try a connection from outside our network to \$webserver in the DMZ on port 25 (SMTP).	The connection should be blocked.
Try a connection from outside our network to \$emailserver in the DMZ on port 80 (HTTP).	The connection should be blocked.
Try a connection from outside our network to \$emailserver in the DMZ on port 25 (SMTP).	The connection should pass.
Try a connection from outside our network to one or more machines in the local network.	The connection should be blocked.

## Debugging Your Rule Set

When your configuration does not behave as expected, there may be an error in the rule set logic, so you need to find the error and correct it. Tracking down logic errors in your rule set can be time-consuming, and could involve manually evaluating your rule set, both as it is stored in the *pf.conf* file and the loaded version after macro expansions and any optimizations.

Users often initially blame PF for problems that turn out to be basic network problems. Network interfaces set to wrong duplex settings, bad netmasks, and faulty network hardware are common culprits.

Before diving into the rule set itself, you can easily determine whether the PF configuration is causing the problem. To do so, disable PF with `pfctl -d` to see if the problem disappears. If the problem persists when PF is disabled, you should turn to debugging other parts of your network configuration instead. On the other hand, if the problem disappeared upon disabling PF, and you are about to start adjusting your PF configuration, make sure that PF is enabled and that your rule set is loaded, with this command:

---

```
$ sudo pfctl -si | grep Status
Status: Enabled for 20 days 06:28:24           Debug: err
```

---

`Status: Enabled` tells us that PF is enabled, so we try viewing the loaded rules with a different `pfctl` command:

---

```
$ sudo pfctl -sr
match in all scrub (no-df max-mss 1440)
```

---

```
block return log all
block return log quick from <bruteforce> to any
anchor "ftp-proxy/*" all
```

---

Here, `pfctl -sr` is equivalent to `pfctl -s` rules. The actual output is likely to be a bit longer than shown here, but this is a good example of what you should expect to see when a rule set is loaded.

For debugging purposes, consider adding the `-vv` flag to the `pfctl` command line to see rule numbers and some additional debug information, like this:

---

```
$ sudo pfctl -vvsr
@0 match in all scrub (no-df max-mss 1440)
  [ Evaluations: 341770  Packets: 3417668  Bytes: 2112276585  States: 125  ]
  [ Inserted: uid 0 pid 14717 State Creations: 92254  ]
@1 match out on nfe0 inet from 10.0.0.0/8 to any queue(q_def, q_pri) nat-to
(nfe0:1) round-robin static-port
  [ Evaluations: 341770  Packets: 0          Bytes: 0          States: 0      ]
  [ Inserted: uid 0 pid 14717 State Creations: 0      ]
@2 match out on nfe0 inet from 192.168.103.0/24 to any queue(q_def, q_pri) nat-to
(nfe0:1) round-robin static-port
  [ Evaluations: 68623   Packets: 2138128   Bytes: 1431276138  States: 103  ]
  [ Inserted: uid 0 pid 14717 State Creations: 39109  ]
@3 block return log all
  [ Evaluations: 341770  Packets: 114929   Bytes: 62705138   States: 0      ]
  [ Inserted: uid 0 pid 14717 State Creations: 0      ]
@4 block return log (all) quick from <bruteforce:>0 to any
  [ Evaluations: 341770  Packets: 2          Bytes: 104        States: 0      ]
  [ Inserted: uid 0 pid 14717 State Creations: 0      ]
@5 anchor "ftp-proxy/*" all
  [ Evaluations: 341768  Packets: 319954   Bytes: 263432399  States: 0      ]
  [ Inserted: uid 0 pid 14717 State Creations: 70     ]
```

---

Now you should perform a structured walk-through of the loaded rule set. Find the rules that match the packets you are investigating. What is the last matching rule? If more than one rule matches, is one of the matching rules a quick rule? (As you probably recall from earlier chapters, when a packet matches a quick rule, evaluation stops, and whatever the quick rule specifies is what happens to the packet.) If so, you will need to trace the evaluation until you hit the end of the rule set or the packet matches a quick rule, which then ends the process. If your rule set walk-through ends somewhere other than the rule you were expecting to match your packet, you have found your logic error.

Rule set logic errors tend to fall into three types:

- Your rule does not match because it is never evaluated. A quick rule earlier in the rule set matched, and the evaluation stopped.
- Your rule is evaluated, but does not match the packet after all, due to the rule's criteria.

- Your rule is evaluated, the rule matches, but the packet also matches another rule later in the rule set. The last matching rule is the one that determines what happens to your connection.

Chapter 8 introduced `tcpdump` as a valuable tool for reading and interpreting PF logs. The program is also very well suited for viewing the traffic that passes on a specific interface. What you learned about PF’s logs and how to use `tcpdump`’s filtering features will come in handy when you want to track down exactly which packets reach which interface.

Here is an example of using `tcpdump` to watch for TCP traffic on the `xlo` interface (but not showing SSH or SMTP traffic) and print the result in very verbose mode (`vvv`):

---

```
$ sudo tcpdump -nvvvi xlo tcp and not port ssh and not port smtp

tcpdump: listening on xlo, link-type EN10MB
21:41:42.395178 194.54.107.19.22418 > 137.217.190.41.80: S [tcp sum ok]
3304153886:3304153886(0) win 16384 <mss 1460,nop,nop,sackOK,nop,wscale
0,nop,nop,timestamp 1308370594 0> (DF) (ttl 63, id 30934, len 64)
21:41:42.424368 137.217.190.41.80 > 194.54.107.19.22418: S [tcp sum ok]
1753576798:1753576798(0) ack 3304153887 win 5792 <mss 1460,sackOK,timestamp
168899231 1308370594,nop,wscale 9> (DF) (ttl 53, id 0, len 60)
```

---

The connection shown here is a successful connection to a website.

There are more interesting things to look for, though, such as connections that fail when they should not according to your specifications, or connections that succeed when your specification says they clearly should not.

The test in these cases involves tracking the packets’ path through your configuration. Once more, it is useful to check to see if PF is enabled or if disabling PF makes a difference. Building on the result from that initial test, you then perform the same kind of analysis of the rule set as described previously:

- Once you have a reasonable theory of how the packets should traverse your rule set and your network interfaces, use `tcpdump` to see the traffic on each of the interfaces in turn.
- Use `tcpdump`’s filtering features to extract only the information you need, to see only the packets that should match your specific case, such as `port smtp` and `dst 192.0.2.19`.
- Find the place where your assumptions no longer match the reality of your network traffic.
- Turn on logging for the rules that may be involved, and turn `tcpdump` loose on the relevant `pflog` interface to see which rule the packets actually match.

The main outline for the test procedure is fairly fixed. If you have narrowed down the cause to your PF configuration, again it’s a case of finding out which rules match and which rule ends up determining whether the packet passes or is blocked.

## **Know Your Network and Stay in Control**

The recurring theme in this book has been how PF and related tools make it relatively easy for you, as the network administrator, to take control of your network and make it behave the way you want it to behave. In other words, build the network you need.

Running a network can be fun, and I hope you have enjoyed this tour of what I consider to be the best tool available for network security. In presenting PF, I have made a conscious decision early on to introduce you to the methods and ways of thinking via interesting and useful configurations, rather than offer a full catalog of available features, or for that matter, make this book the complete reference. The complete PF reference already exists in the man pages, which are updated every six months with the new OpenBSD releases. You can also find further information in the resources I've listed in Appendix A.

Now that you have a broad basic knowledge of what PF can do, you can start building the network according to your own ideas of what you need. You have reached the point where you can find your way around the man pages and locate the exact information you need. This is when the fun part starts!



# A

## RESOURCES



Though I may have wanted to, it proved impossible to cover all possible wrinkles of PF configuration. I hope that the resources listed here will fill in some details or present a slightly different perspective. Some of them are even quite enjoyable reads for their own sake.

### **General Networking and BSD Resources on the Internet**

The following are the general web-accessible resources cited throughout the book. It is worth looking at the various BSD projects' websites for the most up-to-date information.

- Of particular interest for OpenBSD users is the online *OpenBSD Journal* (<http://undeadly.org/>). It offers news and articles about OpenBSD and related issues.

- OpenBSD’s website, <http://www.openbsd.org/>, is the main reference for OpenBSD information. If you’re using OpenBSD, you will be visiting this site every now and then.
- You will find a collection of presentations and papers by OpenBSD developers at <http://www.openbsd.org/papers/>. This site is a good source of information about ongoing developments in OpenBSD.
- *OpenBSD’s Documentation and Frequently Asked Questions*, (<http://www.openbsd.org/faq/index.html>) is more of a user guide than a traditional question-and-answer document. This is where you’ll find a generous helping of background information and step-by-step instructions on how to set up and run your OpenBSD system.
- Henning Brauer’s presentation “Faster Packets—Performance Tuning in the Network Stack and PF” ([http://bulabula.org/papers/2009/eurobsdcon-faster\\_packets/](http://bulabula.org/papers/2009/eurobsdcon-faster_packets/)) is the current main PF developer’s overview of the work done in recent OpenBSD releases to improve network performance, with PF as a main component.
- *PF: The OpenBSD Packet Filter* (<http://www.openbsd.org/faq/pf/index.html>), also known as the *PF User Guide* or the *PF FAQ*, is the official PF documentation, maintained by the OpenBSD team. This guide is updated for each release, and it’s an extremely valuable reference resource for PF practitioners.
- Bob Beck’s “pf. It’s not just for firewalls anymore” (<http://www.ualberta.ca/~beck/nycbug06/pf/>) is a NYCBUG 2006 presentation that covers PF’s redundancy and reliability features, with real-world examples taken from the University of Alberta network.
- Daniel Hartmeier’s PF pages (<http://www.benzedrine.cx/pf.html>) are his collection of PF-related material with links to resources around the Web.
- Daniel Hartmeier’s “Design and Performance of the OpenBSD Stateful Packet Filter (pf)” (<http://www.benzedrine.cx/pf-paper.html>) is the paper he presented at Usenix 2002, which describes the initial design and implementation of PF.
- Daniel Hartmeier’s three-part *undeadly.org* PF series includes “PF: Firewall Ruleset Optimization” (<http://undeadly.org/cgi?action=article&sid=20060927091645>), “PF: Testing Your Firewall” (<http://undeadly.org/cgi?action=article&sid=20060928081238>), and “PF: Firewall Management” (<http://undeadly.org/cgi?action=article&sid=20060929080943>). The three articles cover their respective subjects in great detail, yet manage to be quite readable.
- RFC 1631, The IP Network Address Translator (NAT), May 1994 (<http://www.ietf.org/rfc/rfc1631.txt?number=1631>) is the first part of the NAT specification, which has proved longer lived than the authors had apparently intended. While still an important resource for understanding NAT, it has been largely superseded by the updated RFC 3022, dated January 2001.

- RFC 1918, Address Allocation for Private Internets, February 1996 (<http://www.ietf.org/rfc/rfc1918.txt?number=1918>) is the second part of the NAT and private address space puzzle. This RFC describes the motivations for the allocation of private, nonroutable address space and defines the address ranges. RFC 1918 has been designated a Best Current Practice.
- If you are looking for a text that gives you a thorough and detailed treatment of network protocols with a clear slant toward the TCP/IP worldview, Charles M. Kozierok's *The TCP/IP Guide* (No Starch Press, October 2005), available online with updates at <http://www.tcpipguide.com/>, has few, if any, serious rivals. At more than 1,600 pages, it's not exactly a pocket guide, but it's very useful to have on your desk or in a browser window to set the record straight on any networking terms that you find insufficiently explained in other texts.

## Sample Configurations and Related Musings

A number of people have been kind enough to write up their experiences and make sample configurations available on the Web. The following are some of my favorites.

- Marcus Ranum's "The Six Dumbest Ideas in Computer Security" ([http://www.ranum.com/security/computer\\_security/editorials/dumb/index.html](http://www.ranum.com/security/computer_security/editorials/dumb/index.html)), from September 1, 2005, is a longtime favorite of mine. This article explores some common misconceptions about security and their unfortunate implications for real-world security efforts.
- Randal L. Schwartz's "Monitoring Net Traffic with OpenBSD's Packet Filter" (<http://www.stonehenge.com/merlyn/UnixReview/col51.html>) shows a real-life example of traffic monitoring and using labels for accounting. Some details about PF and labels have changed in the intervening years, but the article is still quite readable and presents several important concepts well.
- The Swedish user group Unix.se's *Brandvägg med OpenBSD* ([http://unix.se/Brandv%4E4gg\\_med\\_OpenBSD](http://unix.se/Brandv%4E4gg_med_OpenBSD)) and its sample configurations such as the basic ALTQ configurations, were quite useful to me early on. The site serves as a nice reminder that volunteer efforts such as local user groups can be excellent sources of information.
- The #pf IRC channel wiki (<http://www.probsd.net/pf/>) is a collection of documentation, sample configurations, and other PF information, maintained by participants in the #pf IRC channel discussions and another example of a very worthwhile volunteer effort.
- Daniele Mazzocchio, an OpenBSD fan from Italy, maintains the website Kernel Panic, with a collection of useful articles and tutorial-like documents on various OpenBSD topics at <http://www.kernel-panic.it/opensbd.html> (in English and Italian). It's well worth the visit for a fresh perspective on various interesting topics from someone who seems to be dedicated to keeping the material up to date with the latest stable OpenBSD versions.

- Randal L. Schwartz’s blog entry for January 29, 2004 (<http://use.perl.org/~merlyn/journal/17094>) shows how he apparently solved an annoying problem via creative use of ALTQ and operating system fingerprinting.
- Kenjiro Cho’s “Managing Traffic with ALTQ” (<http://www.usenix.org/publications/library/proceedings/usenix99/cho.html>) is the original paper that describes the ALTQ design and early implementation on FreeBSD.
- Jason Dixon’s “Failover Firewalls with OpenBSD and CARP,” from *SysAdmin Magazine*, May 2005 (<http://planet.admon.org/howto/failover-firewalls-with-openbsd-and-carp/>) is an overview of CARP and pfsync, with some practical examples.
- Theo de Raadt’s OpenCON 2006 presentation “Open Documentation for Hardware: Why hardware documentation matters so much and why it is so hard to get” (<http://openbsd.org/papers/opencon06-docs/index.html>) was an important inspiration for the note in Appendix B about hardware for free operating systems in general and OpenBSD in particular.

## PF on Other BSD Systems

PF has been ported from OpenBSD to the other BSDs, and while the stated goal for these efforts naturally is to be as up to date as possible in relation to the newest PF versions coming out of OpenBSD, it is useful to keep track of the PF projects in the other BSDs.

- The FreeBSD packet filter (pf) home page (<http://pf4freebsd.love2party.net/>) describes the early work with PF on FreeBSD and the project goals. At the moment, the page is not quite up to date with the latest developments, but it will hopefully spring to life again once Max Laier notices that he is referenced in a printed book.
- The NetBSD project maintains its PF pages at <http://www.netbsd.org/docs/network/pf.html> with updated information about PF on NetBSD.

## BSD and Networking Books

In addition to what appears to be an ever-expanding number of online resources, several books may be useful as companions or supplements to this book.

- Jacek Artymiak, *Building Firewalls with OpenBSD and PF, 2nd Edition* (devGuide.net, 2003). Traditionally the recommended PF book since its first publication, it covers PF in OpenBSD 3.4 in great detail.
- Michael W. Lucas, *Absolute OpenBSD* (No Starch Press, 2003). Written at the time of OpenBSD 3.4, this volume offers a thorough walk-through of OpenBSD with a wealth of hands-on, practical material.
- Michael W. Lucas, *Network Flow Analysis* (No Starch Press, 2010). One of a select few books about network analysis and management using free NetFlow-based tools, this book shows you the tools and methods to discover just what really happens in your network.

- Brandon Palmer and Jose Nazario, *Secure Architectures with OpenBSD* (Addison-Wesley, 2004). This book provides an overview of OpenBSD's features with a marked slant toward building secure and reliable systems. The book references OpenBSD 3.4 as the then up-to-date version.
- Douglas R. Mauro and Kevin J. Schmidt, *Essential SNMP, 2nd Edition* (O'Reilly Media, 2005). As the title says, an essential reference book about SNMP.
- Jeremy C. Reed (editor), *The OpenBSD PF Packet Filter Book* (Reed Media Services, 2006). The book is based on the *PF User Guide*, extended to cover PF on FreeBSD, NetBSD, and DragonFly BSD, and with some additional material on third-party tools that interoperate with PF.
- Christopher M. Buechler and Jim Pingle, *pjSense: The Definitive Guide* (Reed Media Services, 2009). At some 515 pages, this is a comprehensive guide to the FreeBSD and PF-based firewall appliance distribution.

## Wireless Networking Resources

Kjell Jørgen Hole's Wi-Fi courseware (<http://www.kjhole.com/Standards/WiFi/WiFiDownloads.html>) is an excellent resource for understanding wireless networks. The courseware is mainly aimed at University of Bergen students who take Professor Hole's courses, but it is freely available and well worth reading.

For keeping up with developments in the wireless networking world, the WiFiNetNews site offers running updates, with the security-related stories listed at [http://wifinetnews.com/archives/cat\\_security.html](http://wifinetnews.com/archives/cat_security.html).

Another highly recommended resource for wireless security issues is The Unofficial 802.11 Security Web Page (<http://www.drizzle.com/~aboba/IEEE/>).

## spamd and Greylisting-Related Resources

If handling email is part of your life (or is likely to be in the future), you have probably enjoyed the descriptions of `spamd`, tarpitting, and greylisting in this book. If you want a little more background information than what you find in the relevant RFCs, the following documents and web resources provide it.

- Greylisting.org (<http://www.greylisting.org/>) has a useful collection of greylisting-related articles and other information about greylisting and SMTP in general.
- Evan Harris's "The Next Step in the Spam Control War: Greylisting" (<http://greylisting.org/articles/whitepaper.shtml>) is the original greylisting paper.
- Bob Beck's "OpenBSD `spamd`—greylisting and beyond" (<http://www.ualberta.ca/~beck/nycbug06/spamd/>) is an NYCBUG presentation that explains how `spamd` works, leading up to a description of `spamd`'s role in University of Alberta's infrastructure. (Note that much of the "future work" mentioned in the presentation has already been implemented.)

- “Effective spam and malware countermeasures,” <http://home.nuug.no/~peter/foss-aalborg2008/effective-countermeasures.pdf>, originally my BSDCan 2007 paper with some updates, includes a best practice description of how to use greylisting, `spamd`, and various other free tools and OpenBSD to successfully fight spam and malware in your network.

## Book-Related Web Resources

For news and updates about this book, check the book’s home page at the No Starch Press website (<http://www.nostarch.com/pf2.htm>). That page contains links to pages on my personal web space, where various updates and book-related resources will appear as they become available. I will post book-related news and updates at <http://www.bsdly.net/bookofpf/>. Announcements relevant to the book are likely to turn up via my blog at <http://bsdly.blogspot.com/>, too.

I maintain the tutorial manuscript “Firewalling with OpenBSD’s PF packet filter,” which is the forerunner of this book. My policy is to make updates when appropriate, usually as I become aware of changes or features of PF and related software and while preparing for appearances at conferences. The tutorial manuscript is available under a BSD license and can be downloaded in several formats from my web space at <http://home.nuug.no/~peter/pf/>.

Updated versions will appear at that URL more or less in the natural course of tinkering in between events.

## Buy OpenBSD CDs and Donate!

If you have enjoyed this book or found it useful, please go to the OpenBSD.org Ordering page at <http://www.openbsd.org/orders.html> and buy CD sets, or for that matter, support further development work by the OpenBSD project via a donation, such as via the Donations page at <http://www.openbsd.org/donations.html>.

If you are the kind of corporate entity that is more comfortable with donating to a corporation, you can contact the OpenBSD foundation, a Canadian nonprofit corporation created in 2007 for that specific purpose. See the OpenBSD Foundation website at <http://www.openbsdfoundation.org/> for more information.

If you’ve found this book at a conference, there might even be an OpenBSD booth nearby where you can buy CDs, T-shirts, and other items.

Remember that even free software takes real work and real money to develop and maintain.

# B

## A NOTE ON HARDWARE SUPPORT



“How’s the hardware support?” I tend to hear that a lot, and my answer is usually, “In my experience, OpenBSD and other free systems tend to just work.”

But for some reason, there is a general perception that going with free software means that getting hardware components that will actually work will be a serious struggle. In the past, there was some factual basis for this. I remember struggling to install FreeBSD 2.0.5 on the hardware I had available. I was able to boot off the installation CD, but the install never completed because my CD drive was not fully supported.

But that was back in June 1995, when PC CD drives usually came with an almost-but-not-quite IDE interface attached to a sound card, and cheap PCs did not come with networking built in. Configuring a machine for network use usually meant moving jumpers around on the network interface card or the motherboard, or running some weird proprietary setup software—if you had the good luck to be on a system with an Ethernet interface.

Times have changed. Today, you can reasonably expect all important components in your system to work with OpenBSD. Sure, some caution and a bit of planning may be required for building the optimal setup, but that's not necessarily a bad thing.

## Getting the Right Hardware

Getting the right hardware for your system is essentially a matter of checking that your system meets the needs of your project and network:

- Check the online hardware compatibility lists.
- Check the man pages or use `apropos keyword` commands (where `keyword` is the type of device you are looking for).
- Search the archives of relevant mailing lists if you want more background information.
- Use your favorite web search engine to find useful information about how well a specific device works with your operating system.

In most cases, the hardware will work as expected. However, sometimes otherwise functional hardware may come with odd restrictions.

Quite a number of devices are designed to depend on firmware that must be loaded before the operating system can make use of the device. The motivation for this design choice is almost always to lower the cost of the device. When some of those manufacturers refuse to grant redistribution rights for the firmware, the decision becomes a problem because it means that operating systems like OpenBSD can't package the firmware with their releases.

Problems of this type have surfaced in connection with several different types of hardware. In many cases, the manufacturers have been persuaded to change their minds and allow redistribution. However, this does not happen in all cases. One example is the Intel-based wireless networking hardware that is built into many popular laptop models. The hardware is supported in many operating systems, including OpenBSD via the `wpi` and `iwn` drivers. But even with those drivers in place, the hardware will simply not work unless the user has manually fetched and installed the required firmware files.

**NOTE**

*Where supported hardware is restricted, the OpenBSD man pages usually note the fact and may even include the email addresses of people who might be able to change the manufacturer's policy.*

It would take only a minor change in the manufacturer's licensing policy to make life easier for free software users everywhere and boost sales. It's possible that most situations like these will be resolved by the time you read this. Be sure to check the latest information online, and be prepared to vote with your wallet if a particular company refuses to act sensibly.

If you shop online, keep the man pages available in another tab or window. If you go to a physical store, make sure to tell the clerks you will be using a BSD. If you're not sure about the parts they are trying to sell you, ask to borrow a machine to browse the man pages and other documentation

online. You might even ask for permission to boot a machine with the hardware you’re interested in from a CD or USB stick and study the `dmesg` output. Telling shop staff up front about your project could make it easier to get a refund if the part does not work. And if the part does work, letting the vendor know is good advocacy. Your request could very well be the first time they have heard of your favorite operating system.

## Issues Facing Hardware Support Developers

Systems such as OpenBSD and the other BSDs did not spring fully formed from the forehead of a deity (although some will argue that the process was not that different). Rather, they’re the result of years of effort by a number of smart and dedicated developers.

BSD developers are all highly qualified and extremely dedicated people who work tirelessly—the majority in their spare time—to produce amazing results. However, they do not live in a bubble with access to everything they need. The hardware itself or adequate documentation to support it is often unavailable to them. Another common problem is that documentation is often provided only under a nondisclosure agreement (NDA), which limits how developers can use the information.<sup>1</sup>

Through *reverse engineering*, developers can write drivers to support hardware even without proper documentation, but the process is a complicated one that consists of educated guessing, coding, and testing until results begin to emerge. Reverse engineering takes a long time and, for reasons known only to lawmakers and lobbyists, it has legal consequences in several jurisdictions around the world.

The good news is that you can help the developers get the hardware and other material they need.

## How to Help the Hardware Support Efforts

If you can contribute quality code, the BSD projects want to hear from you. If you are not a developer yourself, contributing code may not be an option. Here are several other ways you can contribute:

- *Buy your hardware from open source friendly vendors.* When making decisions or recommendations regarding your organization’s equipment purchases, tell suppliers that *open source friendliness* is a factor in your purchasing decision.
- *Let hardware vendors know what you think about their support (or lack thereof) for your favorite operating system.* Some hardware vendors have been quite helpful, supplying both sample units and programmer documentation. Others have been less forthcoming or downright hostile. Both kinds of vendors, and the ones in between, need encouragement. Write to them

---

1. This is a frequent talk topic too. For example, see Theo de Raadt’s OpenCON 2006 presentation, “Why hardware documentation matters so much and why it is so hard to get,” available at <http://www.openbsd.org/papers/opencon06-docs/index.html>.

and tell them what you think they are doing right and what they can do to improve. If, for example, a vendor has refused to make programming documentation available, or available only under an NDA, a reasoned, well-formulated letter from a potential customer could make the difference.

- *Help test systems and check out the drivers for hardware you are interested in.* If a driver exists or is being developed, the developers are always interested in reports on how their code behaves on other people's equipment. Reports that the system is working fine are always appreciated, but bug reports with detailed descriptions of what goes wrong are even more essential to creating and maintaining a high-quality system.
- *Donate hardware or money.* The developers can always use hardware to develop on, and money certainly helps for day-to-day needs. If you can donate money or hardware, check out the projects' donations page (<http://www.openbsd.org/donations.html> for OpenBSD)<sup>2</sup> or items-needed page (<http://www.openbsd.org/want.html> for OpenBSD). Donations to OpenBSD will most likely help PF development, but if you prefer to donate to FreeBSD, NetBSD, or DragonFly BSD instead, you can find information about how to do so at their websites.

Whatever your relationship with the BSDs and your hardware, I hope that this appendix has helped you to make intelligent decisions about what to buy and how to support the development of the BSDs. Your support will contribute to making more and better quality free software available for everyone.

---

2. If you are a corporate entity that is more comfortable with donating to a corporation, contact the OpenBSD Foundation, a Canadian nonprofit corporation created for that purpose. See the OpenBSD Foundation website at <http://www.openbsdfoundation.org/> for more information.

# INDEX

## Numbers

- 10.0.0.1 IP address, 27
- 192.168.0.1 IP address, 27
- 802.11
  - hardware, 44
  - MAC address filtering, 42–43
  - WEP (Wired Equivalent Privacy), 43
  - WPA (Wi-Fi Protected Access), 43

## Symbols

- < > (angle brackets), 39
- : (colon), 91
- ! (exclamation mark), 39, 58
- ( ) (parentheses), 31, 154, 159

## A

- Absolute OpenBSD* (Lucas), 170
- Acar, Can Erkin, 141
- access points
  - FreeBSD WPA, 48–49
  - with multiple interfaces, 50
  - OpenBSD WPA, 47–48
  - PF rule set, 49–50
- adaptive firewall. *See also* firewall
  - max, 88
  - overload, 88
  - setting up, 86–88
- address-allocation process, 27–28
- addresses
  - IP version 4, 27
  - IP version 6, 27

- nonroutable, 83–84
- routable, 60, 72
- alert syslog level, 156
- (all) logging option, 134–135
- ALTQ (ALTerнате Queueing). *See also* network traffic; traffic
  - ACK packets, 110
  - cbq (class-based queues), 107, 112–113
  - concepts, 106
  - determining bandwidth, 109
  - features of, 105–106
  - on FreeBSD, 107–108
  - handling unwanted traffic, 117–118
  - HFSC (Hierarchical Fair Service Curve)
    - algorithm, 107
    - traffic shaper, 113–115
  - match rules for assignment, 110–111
  - on NetBSD, 108
  - on OpenBSD, 107
  - priq queues, 106–107, 108–111
  - queue concept, 106
  - queue disciplines, 106
  - queue schedulers, 106
  - real-world example, 109–110
  - rule sets, 113
  - for servers in DMZ, 115–117
  - setting up, 107–108
  - syntax for, 108
  - ToS (type of service) fields, 110
  - using to handle traffic, 117–118

ancontrol program, 42  
angle brackets (< >), 39  
antspoof, 159–160  
ARPANET, 27  
ARP balancing, 128  
Artymiak, Jacek, 170  
attack techniques, 159  
`authpf` program, 55–57. *See also*  
    security  
macros and redirection, 57–58  
special-case rules, 57  
    `user_ip` macro, 57  
`auth_web` macro, 58

## B

baseline filtering rule, 19  
Beck, Bob, 104, 168, 171  
Berkeley Software Distribution  
    (BSD) systems, 3, 5  
    vs. Linux, 6  
    reading configurations, 6  
blacklisting mode, setting up  
    `spamd` in, 91–92  
blacklists, 97, 104  
`block in all` rule, 17  
`block-policy` option  
    `drop` value, 152  
    `return` value, 152  
Brauer, Henning, 2, 5, 106,  
    144, 168  
bridge, defined, 78  
bridge setup. *See also* firewall  
    on FreeBSD, 80–81  
    on NetBSD, 81–82  
    on OpenBSD, 79–80  
brute-force attacks, 86–88  
bruteforce table entries,  
    removing, 89  
BSD (Berkeley Software  
    Distribution) systems, 3, 5  
    vs. Linux, 6  
    reading configurations, 6  
Buechler, Christopher M., 171  
*Building Firewalls with OpenBSD and  
    PF, 2nd Edition*  
    (Artymiak), 170  
bytes in and out, showing, 23

## C

CARP (Common Address  
    Redundancy Protocol).  
    *See also* gateways  
`advskew` parameter, 123, 129  
checking configurations, 123  
checking kernel options, 121  
demotion counter, 124  
gateways, 119–121  
    `GENERIC` kernel  
        configurations, 121  
`ifconfig` for interfaces, 122–124  
`ifstated` daemon, 127  
kernel options, 121  
for load balancing, 128–131  
network interfaces, 122–124  
passphrase, 124  
setting `sysctl` values, 121–122  
setting up, 121–124  
`systat` states, 130  
traffic, handling, 126  
using for load balancing,  
    128–130  
`vhid` (virtual host ID), 121  
Cho, Kenjiro, 170  
Christmas Tree EXEC worm, 2  
Cisco’s PIX firewall series  
    exploit, 159  
colon (:), 91  
Common Address Redundancy  
    Protocol (CARP). *See*  
        CARP (Common Address  
            Redundancy Protocol)  
configuration files  
    placement of, 11, 13  
    as program output, 8  
    reading, 6  
configuration tools, 11  
connection information. *See* state  
    table  
content filtering, 90  
control messages. *See* ICMP  
    (Internet Control  
        Message Protocol)  
Core Force firewall product, 5. *See*  
    *also* firewall  
crit syslog level, 156  
`cron` job, creating for `spamd-setup`, 92

## D

debugging. *See also* logging;  
PF (Packet Filter)  
subsystem: logs  
rule sets, 162–164  
using log data for, 150  
debug option, 156–157  
debug syslog level, 156  
deep packet inspection, 2  
Dehmlow, Sven, 2  
demilitarized zone (DMZ), 63–64  
with NAT, 73  
queueing for servers in, 115–117  
testing, 161  
total\_ext bandwidth, 117  
denial-of-service (DoS) attacks,  
83, 159  
de Raadt, Theo, 2, 4, 170  
“Design and Performance of the  
OpenBSD Stateful Packet  
Filter (pf)” (Hartmeier),  
168  
dhclient command, 53  
divert(4) sockets, 2  
Dixon, Jason, 9, 170  
DMZ (demilitarized zone), 63–64  
with NAT, 73  
queueing for servers in, 115–117  
testing, 161  
total\_ext bandwidth, 117  
DNS, running, 60  
domain name resolution,  
performing, 18  
domain names vs. IP addresses, 33  
DoS (denial-of-service) attacks,  
83, 159  
DragonFly BSD, 3, 5

## E

email setups, dealing with, 102–103  
email transmissions, RFC for, 95  
error information, generating,  
156–157  
ESP protocol traffic, 51  
*Essential SNMP, 2nd Edition* (Mauro  
and Schmidt), 171

*/etc/pf.conf* file, 6, 18  
*/etc/rc.conf* file, 6, 13  
*example.com* network, 60  
exclamation mark (!), 39, 58  
“Exploit Mitigation Techniques”  
(de Raadt), 2

## F

failover and redundancy. *See* CARP  
(Common Address  
Redundancy Protocol)  
“Failover Firewalls with OpenBSD  
and CARP” (Dixon), 170  
FAQs (frequently answered  
questions), 7–8  
File Transfer Protocol (FTP), 34  
proxying configuration, 34–36  
security issues, 34  
file transfers, options for, 34  
filtering on interface groups, 76–77  
filtering rules, testing, 23  
firewall, 3. *See also* adaptive firewall;  
bridge setup; Core Force  
firewall product  
configuration  
mistakes in, 26  
to keyword, 26  
guides, applying to rule sets, 21  
implementing as bridges, 78–79.  
Floeter, Reyk, 66  
flowd package  
configuration, 146  
described, 145–146  
filtering features, 147–149  
flows, 146–147  
gateway field, 147  
internalnet macro, 148  
limiting data stored, 148  
protocols, 146  
setting up daemon, 146  
storage of fields in flow  
records, 148  
unwired macro, 148  
verbose output, 147–148  
fragments  
handling, bugs in, 159  
reassembly options, 158

- FreeBSD, 3–4
- ALTQ (ALTernative Queueing)
    - on, 107–108
  - bridge setup on, 80–81
  - connecting to WEP access point, 53–54
  - default values for PF-related settings, 14
  - `/etc/rc.d/pf` script, 15
  - GENERIC kernel, 14
  - IP forwarding, 29
  - kernel options, 121–122
  - packet filter (pf) home page, 170
  - pfSense build, 7
  - rc scripts, 15
  - setting up PF on, 13–15
  - `spamd` in greylisting mode, 96
  - versions of, 14
  - wireless network configuration, 46, 48
- frequently answered questions (FAQs), 7–8
- FTP (File Transfer Protocol), 34
- proxying configuration, 34–36
  - security issues, 34
- ftp-proxy with redirection, 34–36
- ## G
- gateways, 160–161. *See also* CARP (Common Address Redundancy Protocol)
- allowing name service for clients, 32
  - authenticating, setting up, 55–57
  - diagram, 120
  - rule sets, 31
  - setting up, 26, 29–33
  - using pass rule with, 32
- GENERIC kernel, using with FreeBSD, 14
- global settings
- `block-policy`, 152
  - `debug` option, 156–157
  - fragment reassembly, 158
  - `limit` option, 155–156
- optimization option, 158
- reassemble option, 158
- ruleset-optimization option, 157–158
- skip option, 152–153
- state-defaults option, 153–154
- state-policy, 153
- timeout option, 154–155
- greylisting, 93–98
- keeping in sync, 101–102
- mode
- managing, 102–103
  - setting up `spamd` in, 94–96
- web resources, 171
- greytrapping, 98–99, 104
- GUI tool, using with PF rule set, 7
- ## H
- haiku, 8
- hardware support
- developers, 175
  - efforts, 175–176
  - getting, 174–175
- Harris, Evan, 93–94, 171
- Hartmeier, Daniel, 4, 109–110, 141, 168
- Hole, Kjell Jørgen, 42, 171
- hostnames vs. IP addresses, 33
- hosts, providing feedback to, 152
- ## I
- IBM Christmas Tree EXEC worm, 2
- ICMP (Internet Control Message Protocol), 36–37, 39
- codes, 38, 39
  - packet types, 38–39
- IEEE 802.11
- hardware, 44
  - MAC address filtering, 42–43
  - WEP (Wired Equivalent Privacy), 43
  - WPA (Wi-Fi Protected Access), 43
- ifconfig command, 45
- `-a` output of, 30
- bridge configuration, 79–80

CARP configuration, 124  
interface groups, 45, 47, 51, 52,  
  71, 76–77, 123  
ifstated daemon, 127  
ILOVEYOU worm, 2  
in and out rules, 26–27  
information technology (IT), 2  
info syslog level, 156  
inserts counter, 23  
interface groups, filtering on, 76–77  
interface:network notation, 28–29  
interfaces, testing running  
  status of, 30  
interface state daemon, 127  
Internet Control Message Protocol  
  (ICMP), 36–37, 39  
  codes, 38, 39  
  packet types, 38–39  
IP addresses  
  vs. domain names, 33  
  vs. hostnames, 33  
  IPv4 vs IPv6, 27  
  lists of, 39–40  
IPFilter  
  compatibility with, 4  
  copyright infringement  
    episode, 4  
IPSec VPN solutions, 50–51  
iptables vs. PF, 8  
IPv6  
  ICMP updates for, 39  
  vs. NAT (Network Address  
    Translation), 27–28  
  packets, blocking, 23  
  traffic, forwarding, 29  
IT (information technology), 2

## K

KAME project, 28  
keep state, 17  
*kern.debug* log level, 156–157  
kernel  
  hacking, 4  
  memory space, 155  
  PF loadable module, 14, 15, 108  
Knight, Joel, 150  
Kozierok, Charles M., 169

## L

labels, using for traffic statistics,  
  138–139  
Lehey, Greg, 3  
limit option, 155–156  
links, establishing in wireless  
  networks, 42  
Linux  
  vs. BSD, 6  
  naming conventions, 6  
  possibility of running PF on, 7  
lists  
  defined, 18  
  of IP addresses, 39–40  
  maintaining for services, 20  
  managing with *spamdb*, 100–101  
  updating, 101  
  using for readability, 18–22  
load balancing, 66–71, 73–74,  
  128–131  
local network, defining, 28–29  
logging. *See also* debugging; monitoring  
  tools; PF (Packet Filter) subsystem: logs  
  (a11) option, 134–135  
  basics, 132–133  
  data, using for debugging, 150  
  files, rule numbers in, 133  
  legal implications of, 134  
  with *pflogd* daemon, 132  
  to *pflog* interfaces, 135  
  syslog, 135–137  
  using *tcdump* program for, 133  
logical NOT operator (!), 39, 58  
loopback interface, preventing  
  filtering of, 13  
Lucas, Michael W., 170

## M

MAC address filtering, 42–43  
macros  
  using for readability, 18–22,  
    28–29, 31  
  using with *authpf* program,  
    57–58  
mail connections, tracking, 98  
mail-in and mail-out labels, 138

- mail server, 61, 71–72  
 malicious software, 2  
 Management Information Base (MIB), 150  
 “Managing Traffic with ALTQ” (Cho), 170  
 man pages (manuals)  
     consulting, 8  
     looking up, 44  
     listing, 44  
`martians` macro, 83–84  
 match rule, using with `nat-to`, 31  
 Mauro, Douglas R., 171  
 maximum transmission unit (MTU), 38  
 Mazzocchio, Daniele, 169  
 McBride, Ryan, 5  
 memory pools, setting size of, 155–156  
 MIB (Management Information Base), 150  
 Miller, Damien, 2, 145, 150  
 monitoring tools. *See also* logging;  
     NetFlow; PF (Packet Filter) subsystem: logs;  
     `pflow(4)` pseudo-interface  
`flowd`, 145–149  
`flow-tools`, 145  
`nfdump`, 145  
`pfflowd`, 149  
`pftop`, 141  
`pstat`, 141–143  
`systat`, 139–141  
 Morris worm, 2  
 MTU (maximum transmission unit), 38
- N**
- name resolution  
     handling, 20, 60  
     testing, 21–22  
 naming network interfaces, 6  
 NAT (Network Address Translation)  
     DMZ with, 73  
     handling for gateways, 31  
     vs. IPv6, 27–28
- Nazario, Jose, 171  
 NetBSD, 3, 5  
     ALTQ (ALTernative Queueing) on, 108  
 bridge setup on, 81–82  
`/etc/default/pf.boot.conf` file, 16  
`/etc/pf.conf` file, 16  
 IP forwarding, 29  
 kernel options, 121  
 PF pages, 170  
 setting up PF on, 15–16  
 NetFlow, 143–144. *See also*  
     monitoring tools; `pflow(4)`  
     pseudo-interface  
 analysis, 145–149  
 choosing collectors, 145  
 collector and analysis packages, 145  
 data collecting, 145–149, 149–150  
`flowd` package, 145–149  
`flow-tools` package, 145  
`nfdump` package, 145  
 reporting, 145–149  
 sensor, setting up, 144–145  
 network, diagram of, 60, 64, 82, 116, 120, 161  
 Network Address Translation (NAT)  
     DMZ with, 73  
     handling for gateways, 31  
     vs. IPv6, 27–28  
*Network Flow Analysis* (Lucas), 170  
 network interfaces  
     excluding from PF processing, 152–153  
     naming, 6, 31  
 networks  
     with gateways, 120  
     setting up, 74–76  
 network traffic. *See also* ALTQ  
     (ALTernative Queueing); traffic  
     catching via filtering rules, 23  
     cleaning up, 158–160  
     diagram of, 142, 143  
     directing with ALTQ, 105–108  
     IPSec VPN solutions, 50–51

limiting, 3  
logging, 133  
seeing snapshots of, 139–141  
viewing on interfaces, 164  
network troubleshooting, 36–39  
`nixspam` blacklist, 104  
nonroutable addresses, 83–84  
`notice` `syslog` level, 156

## 0

**OpenBSD**  
3.0 base system, 4  
3.1, PF performance, 4  
ALTQ (ALTerate Queueing)  
on, 107  
approach toward design, 2  
approach toward security, 2  
benefits of, 5  
bridge setup on, 79–80  
connecting to WEP access point,  
51–53  
default `pf.conf` file, 13  
encapsulation interface, 51  
`/etc/rc` script, 13  
ifstated daemon, 127  
IP forwarding, 29  
kernel options, 121  
papers by developers, 168  
pass rules, 17  
presentations by developers, 168  
setting up PF on, 12–13  
version of IPFilter, 4  
website, 168  
*OpenBSD Journal*, 167  
optimization option  
aggressive setting, 158  
conservative setting, 158  
high-latency value, 158  
satellite value, 158  
out-of-memory conditions, 159  
out-of-order MX use, detecting, 102

## P

**Packet Filter (PF) subsystem.** *See* PF  
(Packet Filter) subsystem

packets  
displaying live view of, 140  
filtering, 3–4, 18, 76–77  
forwarding, turning on for  
gateways, 29  
getting information about, 23  
logging, 134–135  
matching to state table, 153  
movement, tracking, 137  
normalization, 158–159  
tagging, 77–78  
tracking paths of, 164  
Palmer, Brandon, 171  
parentheses (), 31, 154, 159  
pass rule, using with gateways, 32  
path MTU discovery, 36, 38–39  
Pentium III machine, 5  
permissive rule sets, 19–20  
`pf.conf` file, 13  
`pfctl` program, 11–12  
-d option, 12, 162  
-sm option, 155  
-s timeouts option, 154–155  
using to extract information,  
22–23  
using with tables, 89  
`pfflowd` package, 149–150  
`#pfIRC` channel wiki, 169  
`pflog` interfaces  
cloning, 135  
disabling data accumulation, 136  
logging to, 135  
`pflogd` logging daemon, 132  
`pflow(4)` pseudo-interface, 143–144.  
*See also* monitoring tools;  
NetFlow  
`pflow` device, enabling, 144–145  
`pflow` state option, 143  
PF (Packet Filter) subsystem, 1  
code, finding, 5  
configuration  
converting other products to,  
7–8  
debugging, 162  
confirming running status of, 22  
data, graphing, 141–143  
disabling, 12

PF (Packet Filter) subsystem,  
*continued*  
enabling, 12–13, 162–163  
haiku, 8  
vs. iptables, 8  
logs, 132–133. *See also*  
    debugging; logging;  
    monitoring tools; syslog  
collecting data for, 132  
storage of data, 132  
tracking statistics for rules,  
    137–139  
using labels with, 137–139  
operating system fingerprinting,  
    118  
releases 4.4 through 4.8, 5  
requirements for, 5  
rise of, 3–5  
rules, changes to syntax, 8  
rule set, managing, 7  
running on Linux, 7  
setting up on FreeBSD, 13–15  
setting up on NetBSD, 15–16  
setting up on OpenBSD, 12–13  
user guide, 75  
    version in OpenBSD 4.8, 5  
`pf_rules=` setting, 13  
pfSense build of FreeBSD, 7  
*pfSense: The Definitive Guide*  
    (Buechler and  
    Pingle), 171  
pfstat utility  
    collect statements, 142  
    color values in graphs, 142  
    described, 141  
    home page, 143  
    image definition, 142  
    setting up, 142  
    specifying graph size, 142  
pfsync interfaces, configuring, 125  
pfsync protocol, 119  
    adding, 125–126  
    rule sets, 126–127  
    sysat states, 126  
pftop tool, 141  
ping command, 37  
Pingle, Jim, 171

ping of death, 36  
PIX firewall series exploit, 159  
pool memory, availability of, 155  
PPP connection, using with  
    gateways, 30  
PPPoE, using with gateways, 30  
psstat tool, 141–143  
“Puffy at Work—Getting Code  
    Right and Secure, the  
    OpenBSD Way” (Brauer  
    and Dehmlow), 2

## Q

queues. *See* ALTQ (ALTernate Queueing)  
quick keyword, 32–33

## R

Ranum, Marcus, 2, 18, 169  
readability, using lists and macros  
    for, 18–22  
Realtek Ethernet cards, 31  
reassemble option, 158  
redirection  
    for load balancing, 73–74  
    to pool of addresses, 65–66  
    using with authpf program,  
        57–58  
    using with auth\_web macro, 58  
    using with ftp-proxy, 34–36  
re driver, 26–27  
redundancy and failover. *See* CARP  
    (Common Address  
    Redundancy Protocol)  
Reed, Darren, 4  
Reed, Jeremy C., 171  
relayd daemon  
    CARP-based failover, 71  
    enabling at startup, 69  
    redirects and relays, 66  
    ssl options, 70–71  
    starting, 68  
    sticky-address option, 68  
    tcp options, 71  
    using for load balancing, 128  
webpool table, 68

- remote X11 traffic, blocking, 13  
 removals counter, 23  
 resource exhaustion, 159  
 RFCs  
     114 (FTP), 34  
     765 and 775 (TCP/IP), 34  
     792, 39  
     950, 39  
     1067 (SNMP), 150  
     1191, 39  
     1256, 39  
     1631 (IP NAT), 168  
     1631 (NAT), 28  
     1885 (ICMP updates for IPv6),  
         39  
     1918 (address allocation), 28,  
         60, 169  
     2018, 71  
     2281 (VRRP), 119  
     2460 (IPv6), 28  
     2463 (ICMP updates for IPv6),  
         39  
     2466 (ICMP updates for IPv6),  
         39  
     2521, 39  
     2765, 39  
     2821, 95  
     3330, 60  
     3411 through 3418 (SNMP), 150  
     3768 (VRRP), 119  
     5321, 95  
 Ritschard, Pierre-Yves, 66  
 round-robin option, 65  
 routable addresses, 60, 72  
 rule numbers, displaying for  
     debugging, 163  
 rules  
     changing order of, 157  
     evaluating for gateways, 32  
     expansion of, 138–139  
     getting log data for, 132  
     merging into tables, 157  
     parsing without loading, 21  
     reading, 21  
     removing duplicates, 157  
     removing subsets of, 157  
     tracking statistics for, 137–139
- ruleset-optimization option,  
     157–158  
 rule sets  
     bridge, 82–83  
     building, 16–17  
     checking changes to, 21  
     debugging, 162–164  
     escapes from sequences, 32–33  
     examining, 13  
     firewall considerations, 21  
     keep state part, 17  
     loading, 12, 162–163  
     logic errors, 163–164  
     permissive, 19–20  
     quick keyword, 32–33  
     storage of, 11  
     test case sequence, 162  
     testing, 18  
         after changing, 21–22  
         for gateways, 33  
 Russian name server example, 134

## §

- Schmidt, Kevin J., 171  
 Schwartz, Randal L., 169, 170  
 scrub feature, 158–159  
*Secure Architectures with OpenBSD*  
     (Palmer and Nazario), 171  
 Secure Shell (SSH) service, 86  
 security. *See also* authpf program  
     OpenBSD’s approach to, 2  
     in wireless networks, 42  
“Security Measures in OpenSSH”  
     (Damien Miller), 2  
 Sender Policy Framework (SPF)  
     records, storage of, 103  
 services  
     maintaining lists of, 20  
     running, 65  
     segregating, 63–65  
 set options, 152  
 setup, testing, 160–162  
 Simple Network Management  
     Protocol (SNMP), 150  
 skip option, 152–153

- S**MTP  
 servers, outgoing, 103  
 standards, interpreting, 93–97  
 traffic, initiating, 61–62  
**S**NMP (Simple Network Management Protocol), 150  
 software, malicious, 2  
 spam, fighting, 104  
 SpamAssassin, 90  
 spambd, using to manage lists, 100–101  
**s**pam daemon  
 features of, 89–90  
 keeping greylists in sync, 101–102  
 logging, 93  
 running, 104  
 setting up in blacklisting mode, 91–92  
 setting up in greylisting mode, 94–96  
**s**pamlogd whitelist updater, 98  
**S**PF (Sender Policy Framework)  
 records, storage of, 103  
 spoofing, protecting against, 159–160  
**S**SH brute-force attacks, 86  
**S**SH (Secure Shell) service, 86  
**s**tate-defaults option, 153–154  
**s**tate information, keeping, 17  
**s**tate-policy option  
 floating value, 153  
 if-bound value, 153  
**s**tate table, 17, 153  
 graphing, 142, 143  
 statistics, interpreting, 23  
 viewing, 139–140  
**s**tate-timeout handling, 158. *See also* timeout option  
**s**tate-tracking options, 87–88  
**s**tatistics, displaying live view of, 140  
**s**ticky-address option, 65–66  
**s**tuttering, 90  
**S**YN-flood attacks, 62  
**s**ynproxy state option, 62  
**s**ysctl command, using with IPv6 traffic, 29  
**s**yslog  
 levels, 156  
 logging to, 135–137. *See also* PF (Packet Filter)  
 subsystem: logs  
**s**ystat program, 111  
 bytes view, 140  
 cycling through views, 141  
 iostat view, 141  
 netstat view, 141  
 packets view, 140  
 pf view, 140  
 rules output, 140  
 states output, 139–140  
 vmsstat view, 141  
**s**ystem information, displaying, 22–23  
**s**ystem status. *See* monitoring tools
- T**
- t**ables  
 entries, expiring, 89  
 tidying with pfctl, 89  
 using as lists of IP addresses, 39–40  
**t**ags, 77–78  
**t**arpitting, 90  
**t**cddump program, 133  
 nohup command, 137  
 using to view traffic, 164  
 using with syslog, 136–137  
**T**CPIP  
 configuring client for, 53  
 packet filtering, 30–31, 34, 38, 169  
**T**CP traffic, viewing, 164  
**T**CP vs. UDP services, 20  
**t**esting setups, 160–162  
“**T**he Next Step in the Spam Control War: Greylisting” (Harris), 93–94, 171  
*The OpenBSD PF Packet Filter Book* (Reed), 171  
“**T**he Six Dumbest Ideas in Computer Security” (Ranum), 2, 18, 169  
*The TCP/IP Guide* (Kozierok), 169

timeout option. *See also* state-timeout handling  
adaptive values, 154  
frag value, 154  
inspecting settings for parameters, 154–155  
interval value, 154  
src.track value, 154  
to keyword, with firewalls, 26  
traceroute command, 37–38  
traffic. *See also* ALTQ (ALTerminate Queueing); network traffic  
catching via filtering rules, 23  
cleaning up, 158–160  
diagnostic, permitting, 37  
directing with ALTQ, 105–108  
displaying live view of, 140  
graphing with pfstat, 142, 143  
limiting, 3  
logging, 133  
seeing snapshots of, 139–141  
shaping  
  cbq (class-based queues), 107, 112–113  
  concepts, 106  
  features of, 105–106  
  HFSC (Hierarchical Fair Service Curve), 107  
  queue concept, 106  
  queue disciplines, 106  
  queue schedulers, 106  
  real-world example, 109–110  
  setting up, 107–108  
  ToS (type of service)  
    fields, 110  
  using to handle traffic, 117–118  
showing snapshots of, 141  
totals, 137  
viewing on interfaces, 164  
traplist, setting up, 99–100  
trojans, 2  
troubleshooting networks  
  ICMP protocol, 36–37  
  path MTU discovery, 38–39  
  ping command, 37  
  traceroute command, 37–38

## U

UDP vs. TCP services, 20  
Unix.se user group, 169  
*/usr/share/examples/pf/pf.conf* file, 15

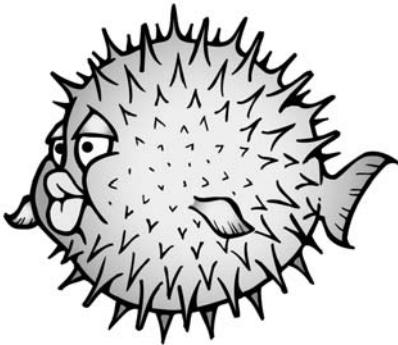
## V

verbose mode, 20  
virtual local area network (VLAN), 63  
virtual private networks (VPNs), setting up, 50–51  
Virtual Router Redundancy Protocol (VRRP), 119  
viruses, 2  
VLAN (virtual local area network), 63  
VPNs (virtual private networks), setting up, 50–51  
VRRP (Virtual Router Redundancy Protocol), 119

## W

warning syslog level, 156  
webpool table, creating, 68  
web server, running, 71–72  
websites  
  Cisco’s PIX firewall series  
    exploit, 159  
  “Explaining BSD,” 3  
  flow-tools package, 145  
  FreeBSD packet filter (pf) home page, 170  
  greylisting.org, 171  
  Hartmeier, Daniel, 4  
  network security, 42  
  nfdump package, 145  
  OpenBSD, 168  
  OpenBSD security, 2  
  pfSense (FreeBSD build), 7  
  security, 42  
  SpamAssassin, 90  
  Wi-Fi Net News, 42  
WEP (Wired Equivalent Privacy), 43, 45  
wicontrol program, 42  
Wi-Fi Net News website, 42

- Wi-Fi Protected Access (WPA), 43, 47–48
- Wired Equivalent Privacy (WEP), 43, 45
- wireless networks
  - access points
    - FreeBSD WPA, 48–49
    - with multiple interfaces, 50
    - OpenBSD WPA, 47–48
    - PF rule set, 49–50
  - client side, 51
  - establishing links in, 42
  - FreeBSD WEP setup, 46
- FreeBSD WPA access point, 48–49
- IPSec VPN solutions, 50–51
- OpenBSD WEP setup, 44
- OpenBSD WPA access point, 47–48
- security in, 42
- setting up, 44–46
- viewing kernel messages, 44
- worms, 2
- `wpa-psk` utility, running, 47
- `wpa_supplicant`, setting up, 54
- WPA (Wi-Fi Protected Access), 43, 47–48



# THE OPENBSD FOUNDATION

## A CANADIAN NOT-FOR-PROFIT CORPORATION

OPENBSD · OPENSSH · OPENBGPD · OPENNTPD · OPENCVS

The OpenBSD Foundation exists to support OpenBSD—the home of pf—and related projects. While the OpenBSD Foundation works in close cooperation with the developers of these wonderful free software projects, it is a separate entity.

If you use pf in a corporate environment, please point management to the URL below, and encourage them to contribute financially to the Foundation.

[www.OPENBSDFOUNDATION.ORG](http://www.OPENBSDFOUNDATION.ORG)





**The Electronic Frontier Foundation (EFF)** is the leading organization defending civil liberties in the digital world. We defend free speech on the Internet, fight illegal surveillance, promote the rights of innovators to develop new digital technologies, and work to ensure that the rights and freedoms we enjoy are enhanced — rather than eroded — as our use of technology grows.

## PRIVACY

EFF has sued telecom giant AT&T for giving the NSA unfettered access to the private communications of millions of their customers. [eff.org/nsa](http://eff.org/nsa)

## FREE SPEECH

EFF's Coders' Rights Project is defending the rights of programmers and security researchers to publish their findings without fear of legal challenges. [eff.org/freespeech](http://eff.org/freespeech)

## INNOVATION

EFF's Patent Busting Project challenges overbroad patents that threaten technological innovation. [eff.org/patent](http://eff.org/patent)

## FAIR USE

EFF is fighting prohibitive standards that would take away your right to receive and use over-the-air television broadcasts any way you choose. [eff.org/IP/fairuse](http://eff.org/IP/fairuse)

## TRANSPARENCY

EFF has developed the Switzerland Network Testing Tool to give individuals the tools to test for covert traffic filtering. [eff.org/transparency](http://eff.org/transparency)

## INTERNATIONAL

EFF is working to ensure that international treaties do not restrict our free speech, privacy or digital consumer rights. [eff.org/global](http://eff.org/global)

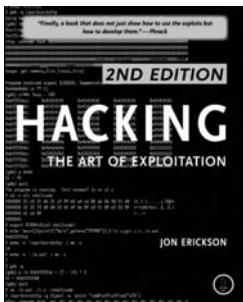
EFF.ORG

ELECTRONIC FRONTIER FOUNDATION

Protecting Rights and Promoting Freedom on the Electronic Frontier

EFF is a member-supported organization. Join Now! [www.eff.org/support](http://www.eff.org/support)

More No-Nonsense Books from  NO STARCH PRESS



## HACKING, 2ND EDITION

### The Art of Exploitation

by JON ERICKSON

Whereas many security books merely show how to run existing exploits, *Hacking: The Art of Exploitation* was the first book to explain how exploits actually work—and how you can develop and implement your own. In this all new second edition, author Jon Erickson uses practical examples to illustrate the fundamentals of serious hacking. You'll learn about key concepts underlying common exploits, such as programming errors, assembly language, networking, shellcode, cryptography, and more. And the bundled Linux LiveCD provides an easy-to-use, hands-on learning environment. This edition has been extensively updated and expanded, including a new introduction to the complex, low-level workings of computers.

FEBRUARY 2008, 488 PP. W/CD, \$49.95

ISBN 978-1-59327-144-2



## GRAY HAT PYTHON

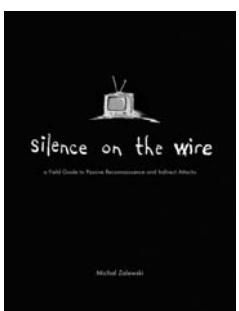
### Python Programming for Hackers and Reverse Engineers

by JUSTIN SEITZ

*Gray Hat Python* explains how to complete various hacking tasks with Python, which is fast becoming the programming language of choice for hackers, reverse engineers, and software testers. Author Justin Seitz explains the concepts behind hacking tools like debuggers, Trojans, fuzzers, and emulators. He then goes on to explain how to harness existing Python-based security tools and build new ones when the pre-built ones just won't cut it. The book teaches readers how to automate tedious reversing and security tasks; sniff secure traffic out of an encrypted web browser session; use PyDBG, Immunity Debugger, Sulley, IDAPython, and PyEMU; and more.

APRIL 2009, 216 PP., \$39.95

ISBN 978-1-59327-192-3



## SILENCE ON THE WIRE

### A Field Guide to Passive Reconnaissance and Indirect Attacks

by MICHAL ZALEWSKI

*Silence on the Wire: A Field Guide to Passive Reconnaissance and Indirect Attacks* explains how computers and networks work, how information is processed and delivered, and what security threats lurk in the shadows. No humdrum technical white paper or how-to manual for protecting one's network, this book is a fascinating narrative that explores a variety of unique, uncommon, and often quite elegant security challenges that defy classification and eschew the traditional attacker-victim model.

APRIL 2005, 312 PP., \$39.95

ISBN 978-1-59327-046-9

## **ABSOLUTE FREEBSD, 2ND EDITION**

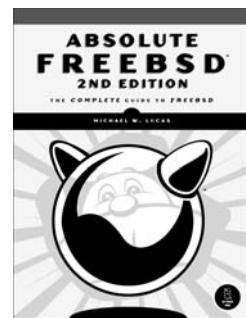
**The Complete Guide to FreeBSD**

by MICHAEL W. LUCAS

*Absolute FreeBSD, 2nd Edition* is the newly updated edition of the best-selling and highly regarded guide to FreeBSD, now covering version 7.0. Written by FreeBSD committer Michael W. Lucas with the help and advice of dozens of FreeBSD developers, *Absolute FreeBSD, 2nd Edition* covers installation, networking, security, network services, system performance, kernel tweaking, filesystems, SMP, upgrading, crash debugging, and much more.

NOVEMBER 2007, 744 PP., \$59.95

ISBN 978-1-59327-151-0



## **BUILDING A SERVER WITH FREEBSD 7**

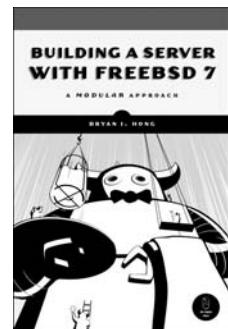
**A Modular Approach**

by BRYAN J. HONG

The most difficult part of building a server with FreeBSD, the Unix-like operating system, is arguably software installation and configuration. Finding the software is easy enough; getting everything up and running is another thing entirely. If you're a small business owner looking for a reliable email server, a curious Windows administrator, or if you just want to put that old computer in the closet to work, *Building a Server with FreeBSD 7* will show you how to get things up and running quickly. You'll learn how to install FreeBSD and then how to install popular server applications with the ports collection. Each package is treated as an independent module, so you can dip into the book at any point to install just the packages you need, when you need them.

APRIL 2008, 288 PP., \$34.95

ISBN 978-1-59327-145-9



**PHONE:**

800.420.7240 OR

415.863.9900

MONDAY THROUGH FRIDAY,

9 A.M. TO 5 P.M. (PST)

**EMAIL:**

SALES@NOSTARCH.COM

**WEB:**

WWW.NOSTARCH.COM

**FAX:**

415.863.9950

24 HOURS A DAY,

7 DAYS A WEEK

**MAIL:**

NO STARCH PRESS

38 RINGOLD STREET

SAN FRANCISCO, CA 94103

USA



*The Book of PF, 2nd Edition* is set in New Baskerville, TheSansMono Condensed, Futura, and Dogma.

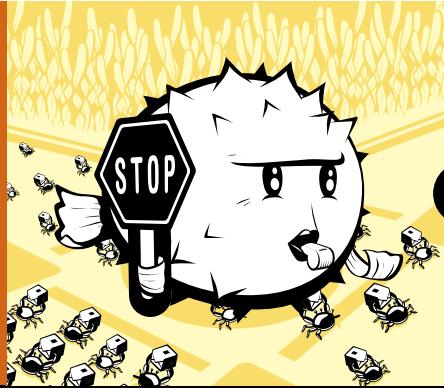
This book was printed and bound by Transcontinental, Inc. at Transcontinental Gagné in Louiseville, Quebec, Canada. The paper is Domtar Husky 60# Smooth, which is certified by the Forest Stewardship Council (FSC). The book has an Otabind binding, which allows it to lie flat when open.

## **UPDATES**

Visit <http://www.nostarch.com/pf2.htm> for updates, errata, and other information.



# BUILD A MORE SECURE NETWORK WITH PF



2ND EDITION

Covers OpenBSD 4.8,  
FreeBSD 8.1, and  
NetBSD 5

OpenBSD's stateful packet filter, PF, is the heart of the OpenBSD firewall and a necessity for any admin working in a BSD environment. With a little effort and this book, you'll gain the insight needed to unlock PF's full potential.

This second edition of *The Book of PF* has been completely updated and revised. Based on Peter N.M. Hansteen's popular PF website and conference tutorials, this no-nonsense guide covers NAT and redirection, wireless networking, spam fighting, failover provisioning, logging, and more. Throughout the book, Hansteen emphasizes the importance of staying in control with a written network specification, keeping rule sets readable using macros, and performing rigid testing when loading new rules.

*The Book of PF* tackles a broad range of topics that will stimulate your mind and pad your resume, including how to:

- Create rule sets for all kinds of network traffic, whether it's crossing a simple LAN, hiding behind NAT, traversing DMZs, or spanning bridges or wider networks
- Create wireless networks with access points, and lock them down with authpf and special access restrictions

- Maximize flexibility and service availability via CARP, relayd, and redirection
- Create adaptive firewalls to proactively defend against would-be attackers and spammers
- Implement traffic shaping and queues with ALTQ (priq, cbq, or hfsc) to keep your network responsive
- Master your logs with monitoring and visualization tools (including NetFlow)

*The Book of PF* is for BSD enthusiasts and network administrators at any skill level. With more and more services placing high demands on bandwidth and an increasingly hostile Internet environment, you can't afford to be without PF expertise.

## ABOUT THE AUTHOR

Peter N.M. Hansteen is a consultant, writer, and sysadmin based in Bergen, Norway. A longtime Freenix advocate, Hansteen is a frequent lecturer on OpenBSD and FreeBSD topics, an occasional contributor to *BSD Magazine*, and one of the original members of the RFC 1149 implementation team. He writes a frequently slashdotted blog (<http://bsdly.blogspot.com/>) and is the author of the highly regarded PF tutorial (<http://home.nuug.no/~peter/pf/>).



THE FINEST IN GEEK ENTERTAINMENT™

[www.nostarch.com](http://www.nostarch.com)



This book uses a lay-flat binding that won't snap shut.

ISBN: 978-1-59327-274-6



9 781593 272746

\$29.95 (\$34.95 CDN)



6 89145 72740 1

SHIPPING:  
OPERATING SYSTEMS/UNIX