

# Network Traffic Analysis with SiLK

Analyst's Handbook for SiLK Version 3.12.0 and Later

**AUGUST 2018**

Paul Krystosek  
Nancy Ott  
Geoffrey Sanders  
Timothy Shimeall

CERT® Situational Awareness Group



**Carnegie Mellon University**  
Software Engineering Institute

**Network Traffic Analysis  
with SiLK**

**ANALYSTS' HANDBOOK**  
for SiLK Versions 3.12.0 and Later

Paul Krystosek  
Nancy M. Ott  
Geoffrey Sanders  
Timothy Shimeall

*August 2018*

**CERT® Situational Awareness Group**

Copyright 2018 Carnegie Mellon University. All Rights Reserved.

This material is based upon work funded and supported by the Department of Homeland Security under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center sponsored by the United States Department of Defense.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

References herein to any specific commercial product, process, or service by trade name, trade mark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by Carnegie Mellon University or its Software Engineering Institute.

This report was prepared for the SEI Administrative Agent AFLCMC/AZS 5 Eglin Street Hanscom AFB, MA 01731-2100

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

Internal use:\* Permission to reproduce this material and to prepare derivative works from this material for internal use is granted, provided the copyright and "No Warranty" statements are included with all reproductions and derivative works.

External use:\* This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other external and/or commercial use. Requests for permission should be directed to the Software Engineering Institute at [permission@sei.cmu.edu](mailto:permission@sei.cmu.edu).

\* These restrictions do not apply to U.S. government entities.

Carnegie Mellon<sup>®</sup>, CERT<sup>®</sup>, CERT Coordination Center<sup>®</sup> and FloCon<sup>®</sup> are registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM18-0954

[DISTRIBUTION STATEMENT A]  
This material has been approved for public  
release and unlimited distribution.

Adobe is a registered trademark of Adobe Systems Incorporated in the United States and/or other countries.  
Akamai is a registered trademark of Akamai Technologies, Inc.  
Apple and OS X are trademarks of Apple Inc., registered in the U.S. and other countries.  
Cisco Systems is a registered trademark of Cisco Systems, Inc. and/or its affiliates in the United States and certain other countries.  
DOCSIS is a registered trademark of CableLabs.  
FreeBSD is a registered trademark of the FreeBSD Foundation.  
IEEE is a registered trademark of The Institute of Electrical and Electronics Engineers, Inc.  
JABBER is a registered trademark and its use is licensed through the XMPP Standards Foundation.  
Linux is the registered trademark of Linus Torvalds in the U.S. and other countries.  
MaxMind, GeoIP, GeoLite, and related trademarks are the trademarks of MaxMind, Inc.  
Microsoft and Windows are registered trademarks of Microsoft Corporation in the United States and/or other countries.  
OpenVPN is a registered trademark of OpenVPN Technologies, Inc.  
Perl is a registered trademark of The Perl Foundation.  
Python is a registered trademark of the Python Software Foundation.  
Snort is a registered trademark of Cisco and/or its affiliates.  
Solaris is a registered trademark of Oracle and/or its affiliates in the United States and other countries.  
UNIX is a registered trademark of The Open Group.  
VPNz is a registered trademark of Advanced Network Solutions, Inc.  
Wireshark is a registered trademark of the Wireshark Foundation.  
All other trademarks are the property of their respective owners.



# Contents

<b>Contents</b>	<b>v</b>
<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>List of Examples</b>	<b>xv</b>
<b>Acknowledgements</b>	<b>xix</b>
<b>Handbook Goals</b>	<b>xxi</b>
<b>1 Introduction to SiLK</b>	<b>1</b>
1.1 What is SiLK?	1
1.2 The SiLK Flow Repository	2
1.2.1 What is Network Flow Data?	2
1.2.2 Structure of a Flow Record	3
1.2.3 Flow Generation and Collection	3
1.2.4 Introduction to Flow Collection	6
1.2.5 Where Network Flow Data Are Collected	6
1.2.6 Types of Network Traffic	7
1.2.7 The Collection System and Data Management	7
1.2.8 How Network Flow Data Are Organized	8
1.3 The SiLK Tool Suite	8
1.4 How to Use SiLK for Analysis	9
1.4.1 Single-path Analysis	9
1.4.2 Multi-path Analysis	9
1.4.3 Exploratory Analysis	10
1.5 Workflow for SiLK Analysis	10
1.5.1 Formulate	10
1.5.2 Model	11
1.5.3 Test	12
1.5.4 Analyze	12
1.5.5 Refine	12
1.6 Applying the SiLK Workflow	12
1.7 Dataset for Single-path, Multi-path, and Exploratory Analysis Examples	13

<b>2</b>	<b>Basic Single-path Analysis with SiLK: Profiling and Reacting</b>	<b>15</b>
2.1	Single-path Analysis: Concepts	15
2.1.1	Scoping Queries of Network Flow Data	16
2.1.2	Excluding Unwanted Network Traffic	17
2.1.3	Example Single-Path Analysis	17
2.2	Single-path Analysis: Analytics	17
2.2.1	Get a List of Sensors With <code>rwsiteinfo</code>	17
2.2.2	Choose Flow Records With <code>rwfilter</code>	21
2.2.3	View Flow Records With <code>rwcut</code>	24
2.2.4	Viewing File Information with <code>rwfileinfo</code>	26
2.2.5	Profile Flows With <code>rwuniq</code> and <code>rwstats</code>	28
2.2.6	Characterize Traffic by Time Period With <code>rwcount</code>	32
2.2.7	Sort Flow Records With <code>rwsort</code>	35
2.2.8	Use IPsets to Gather IP Addresses	37
2.2.9	Resolve IP Addresses to Domain Names With <code>rwresolve</code>	40
<b>3</b>	<b>Case Studies: Basic Single-path Analysis</b>	<b>43</b>
3.1	Profile Traffic Around an Event	43
3.1.1	Examining Shifts in Traffic	44
3.1.2	How to Profile Traffic	45
3.2	Generate Top <i>N</i> Lists	46
3.2.1	Using <code>rwstats</code> to Create Top <i>N</i> Lists	46
3.2.2	Interpreting the Top- <i>N</i> Lists	49
<b>4</b>	<b>Intermediate Multi-path Analysis with SiLK: Explaining and Investigating</b>	<b>51</b>
4.1	Multi-path Analysis: Concepts	51
4.1.1	What Is Multi-path Analysis?	51
4.1.2	Example of a Multi-path Analysis: Examining Web Service Traffic	52
4.1.3	Exploring Relationships and Behaviors With Multi-path Analysis	53
4.1.4	Integrating and Interpreting the Results of Multi-path Analysis	55
4.1.5	“Gotchas” for Multi-path Analysis	55
4.2	Multi-path Analysis: Analytics	56
4.2.1	Complex Filtering With <code>rwfilter</code>	56
4.2.2	Finding Low-Packet Flows with <code>rwfilter</code>	61
4.2.3	Time Binning, Options, and Thresholds With <code>rwstats</code> , <code>rwuniq</code> and <code>rwcount</code>	63
4.2.4	Summarizing Network Traffic with Bags	67
4.2.5	Working with Bags and IPsets	74
4.2.6	Masking IP Addresses	75
4.2.7	Working With IPsets	76
4.2.8	Indicating Flow Relationships	82
4.2.9	Managing IPset, Bag, and Prefix Map Files	90
<b>5</b>	<b>Case Studies: Intermediate Multi-path Analysis</b>	<b>93</b>
5.1	Building Inventories of Network Flow Sensors With IPsets	93
5.1.1	Path 1: Associate Addresses with a Single Sensor	94
5.1.2	Path 2: Associate Addresses of Remaining Sensors	95
5.1.3	Path 3: Associate Shared Addresses	95
5.1.4	Merge Address Results	96
5.2	Automating IPset Inventories of Network Flow Sensors	96
5.2.1	Program Header	96



5.2.2	Program Loop	97
<b>6</b>	<b>Advanced Exploratory Analysis with SiLK: Exploring and Hunting</b>	<b>99</b>
6.1	Exploratory Analysis: Concepts	99
6.1.1	Exploring Network Behavior	100
6.1.2	Starting Points for Exploratory Analysis	101
6.1.3	Example Exploratory Analysis: Investigating Anomalous NTP Activity	101
6.1.4	Observations on Exploratory Analysis	106
6.2	Exploratory Analysis: Analytics	107
6.2.1	Using Tuple Files for Complex Filtering	107
6.2.2	Manipulating Bags	108
6.2.3	Sets Versus Bags: A Scanning Example	112
6.2.4	Manipulating SiLK Files	114
6.2.5	Dividing or Sampling Flow Record Files with <code>rwsplit</code>	116
6.2.6	Generate Flow Records From Text	119
6.2.7	Labeling Data with Prefix Maps	121
6.2.8	Translating IDS Signatures into <code>rwfilter</code> Calls	130
<b>7</b>	<b>Case Studies: Advanced Exploratory Analysis</b>	<b>131</b>
7.1	Level 0: Which TCP Requests are Suspicious?	133
7.2	Level 1: How Can We Identify and React to Illegitimate Requests?	135
7.3	Level 2: What are the Illegitimate Sources and Destinations Doing?	137
7.3.1	Level 2A: What are the Illegitimate Source IPs Doing?	137
7.3.2	Level 2B: What Behavior Changes do Destination IPs Show?	138
7.4	Level 3: What are the Commonalities Across The Cases?	140
<b>8</b>	<b>Extending the Reach of SiLK with PySiLK</b>	<b>141</b>
8.1	Using PySiLK	142
8.1.1	PySiLK Requirements	142
8.1.2	PySiLK Scripts and Plug-ins	142
8.2	Extending <code>rwfilter</code> with PySiLK	143
8.2.1	Using PySiLK to Incorporate State from Previous Records: Eliminating Inconsistent Sources	143
8.2.2	Using PySiLK to Incorporate State from Previous Records: Detecting Port Knocking	145
8.2.3	Using PySiLK with <code>rwfilter</code> in a Distributed or Multiprocessing Environment	147
8.2.4	Simple PySiLK with <code>rwfilter --python-expr</code>	147
8.2.5	PySiLK with Complex Combinations of Rules	148
8.2.6	Use of Data Structures in Partitioning	148
8.3	Extending SiLK with Fields Defined with PySiLK	151
8.4	Extending <code>rwcut</code> and <code>rwsort</code> with PySiLK	151
8.4.1	Computing Values from Multiple Records	152
8.4.2	Computing a Value Based on Multiple Fields in a Record	152
8.4.3	Defining a Character String Field for <code>rwcut</code>	154
8.4.4	Defining a Character String Field for Five SiLK Tools	154
8.5	Defining Key Fields and Summary Value Fields for <code>rwuniq</code> and <code>rwstats</code>	158

<b>A</b>	<b>Networking Primer</b>	<b>161</b>
A.1	Understanding TCP/IP Network Traffic	161
A.2	TCP/IP Protocol Layers	161
A.3	Structure of the IP Header	163
A.4	IP Addressing and Routing	164
A.4.1	Structure of an IP Address	164
A.4.2	Reserved IP Addresses	165
A.5	Major Protocols	168
A.5.1	Protocol Layers and Encapsulation	168
A.5.2	Transmission Control Protocol (TCP)	168
A.5.3	UDP and ICMP	171
<b>B</b>	<b>Using UNIX to Implement Network Traffic Analysis</b>	<b>173</b>
B.1	Using the UNIX Command Line	173
B.2	Standard In, Out, and Error	176
B.2.1	Output Redirection	176
B.2.2	Input Redirection	176
B.2.3	Pipes	177
B.2.4	Here-Documents	177
B.2.5	Named Pipes	178
B.3	Script Control Structures	179
<b>C</b>	<b>SiLK Commands</b>	<b>181</b>
C.1	Getting Help with SiLK Tools	181
C.2	rwsiteinfo Command Summary	182
C.3	rwfilter Command Summary	183
C.4	rwstats Command Summary	190
C.5	rwcount Command Summary	191
C.6	rwcut Command Summary	192
C.7	rwsort Command Summary	194
C.8	rwuniq Command Summary	195
C.9	rwnetmask Command Summary	196
C.10	rwcat Command Summary	197
C.11	rwappend Command Summary	198
C.12	rwsplit Command Summary	199
C.13	rwputc Command Summary	200
C.14	rwset Command Summary	201
C.15	rwsetcat Command Summary	202
C.16	rwsettool Command Summary	203
C.17	rwsetbuild Command Summary	204
C.18	rwbag Command Summary	205
C.19	rwbagbuild Command Summary	206
C.20	rwbagcat Command Summary	208
C.21	rwbagtool Command Summary	209
C.22	rwfileinfo Command Summary	210
C.23	rwmapbuild Command Summary	211
C.24	rwmaplookup Command Summary	212
C.25	rwmatch Command Summary	213
C.26	rwgroup Command Summary	214
C.27	Features Common to Several Commands	215

C.27.1 Parameters Common to Several Commands . . . . .	215
<b>D Additional Information on SiLK</b>	<b>219</b>
D.1 SiLK Support and Documentation . . . . .	219
D.2 FloCon Conference and Social Media . . . . .	220
D.3 Email Addresses and Mailing Lists . . . . .	220
<b>E Further Reading and Resources</b>	<b>223</b>
E.1 Network Flow and Related Topics . . . . .	223
E.1.1 Technical Papers . . . . .	223
E.1.2 Books on Network Flow and Network Security . . . . .	224
E.2 Bash Scripting Resources . . . . .	224
E.2.1 Online Tutorial . . . . .	224
E.2.2 Books on Bash Scripting . . . . .	224
E.3 Visualization . . . . .	225
E.3.1 Rayon . . . . .	225
E.3.2 FloViz . . . . .	225
E.3.3 Graphviz - Graph Visualization Software . . . . .	225
E.3.4 The Spinning Cube of Potential Doom . . . . .	225
E.4 Networking Standards . . . . .	226
<b>Index</b>	<b>227</b>



# List of Figures

1.1	From Packets to Flows . . . . .	5
1.2	Default Traffic Types for Sensors . . . . .	6
1.3	SiLK Analysis Workflow . . . . .	11
2.1	Single-Path Analysis . . . . .	16
2.2	<code>rwfilter</code> Parameter Relationships . . . . .	22
2.3	Displaying <code>rwcount</code> Output Using 10-Minute and 1-Minute Bins . . . . .	34
4.1	Multi-Path Analysis . . . . .	52
4.2	Diagram of a Simple, Non-overlapping Manifold . . . . .	57
4.3	Diagram of a Complex, Overlapping Manifold . . . . .	57
4.4	Client and Server TCP flags . . . . .	59
4.5	Allocating Flows, Packets and Bytes via <code>rwcount</code> Load-Schemes . . . . .	63
6.1	Exploratory Analysis . . . . .	100
6.2	Time Series Plot of NTP Traffic . . . . .	104
7.1	FCC Network Diagram . . . . .	132
A.1	TCP/IP Protocol Layers . . . . .	162
A.2	Structure of the IPv4 Header . . . . .	163
A.3	TCP Header . . . . .	168
A.4	TCP State Machine . . . . .	170
A.5	UDP and ICMP Headers . . . . .	171
C.1	<code>rwfilter</code> Partitioning Parameters . . . . .	184



# List of Tables

1.1	Fields in a SiLK Network Flow record	4
A.1	IPv4 Reserved Addresses	166
A.2	IPv6 Reserved Addresses	167
B.1	Some Common UNIX Commands	174
C.1	Parameters for <code>rwsiteinfo --fields</code>	182
C.2	<code>rwfilter</code> Selection Parameters	183
C.3	Single-Integer- or Range-Partitioning Parameters	184
C.4	Multiple-Integer- or Range-Partitioning Parameters	185
C.5	Address-Partitioning Parameters	185
C.6	High/Mask Partitioning Parameters	185
C.7	Time-Partitioning Parameters	186
C.8	Prefix-Map-Partitioning Parameters	186
C.9	Miscellaneous Partitioning Parameters	186
C.10	<code>rwfilter</code> Output Parameters	188
C.11	Miscellaneous <code>rwfilter</code> Parameters	189
C.12	Time distribution options for <code>rwcount --load-scheme</code>	191
C.13	Arguments for the <code>--fields</code> Parameter	193
C.14	Output-Filtering Options for <code>rwuniq</code>	195
C.15	Fixed-Value Parameters for <code>rwtuc</code>	200
C.16	<code>rwbagbuild</code> Key or Value Options	207
C.17	Common Parameters in Essential SiLK Tools	216
C.18	Parameters Common to Several Commands	217
C.19	<code>--ip-format</code> Values	218
C.20	<code>--timestamp-format</code> <i>format</i> , <i>modifier</i> , and <i>timezone</i> Values	218





# List of Examples

2.1	Using <code>rwsiteinfo</code> to List Sensors, Display Traffic Types, and Show Repository Information .	19
2.2	Using <code>rwfilter</code> to Retrieve Network Flow Records From The SiLK Repository . . . . .	23
2.3	<code>rwcut</code> for Displaying the Contents of Ten Flow Records . . . . .	25
2.4	<code>rwcut --fields</code> to Rearrange Output . . . . .	26
2.5	<code>rwfileinfo</code> Displays Flow Record File Characteristics . . . . .	27
2.6	Characterizing flow byte counts with <code>rwuniq</code> . . . . .	30
2.7	Finding the top protocols with <code>rwstats</code> . . . . .	31
2.8	Counting Bytes, Packets and Flows with Respect to Time . . . . .	33
2.9	Sorting by Destination IP Address, Protocol, and Byte Count . . . . .	36
2.10	Using <code>rwset</code> to Gather IP Addresses . . . . .	37
2.11	Using <code>rwsetbuild</code> to Gather IP Addresses . . . . .	38
2.12	Using <code>rwsetcat</code> to Count Gathered IP Addresses . . . . .	39
2.13	Using <code>rwsetcat</code> to Print Networks and Host Counts . . . . .	39
2.14	Using <code>rwsetcat</code> to Print IP Address Statistical Summaries . . . . .	40
3.1	Using <code>rwfilter</code> and <code>rwuniq</code> to Profile Traffic Around an Event . . . . .	45
3.2	Collated Profile of Traffic Around an Event . . . . .	46
3.3	Removing Unneeded Flows for Top <i>N</i> . . . . .	48
4.1	Examining Flows for Web Service Ports . . . . .	54
4.2	Simple Manifold to Select Inbound Client and Server Flows . . . . .	58
4.3	Complex Manifold to Select Inbound Client and Server Flows . . . . .	60
4.4	Extracting Low-Packet Flow Records . . . . .	62
4.5	Constraining Counts to a Threshold by using <code>rwuniq --flows</code> . . . . .	64
4.7	Constraining Flow and Packet Counts with <code>rwuniq --flows</code> and <code>--packets</code> . . . . .	64
4.6	Setting Minimum Flow Thresholds with <code>rwuniq --values</code> . . . . .	65
4.8	Profiling IP addresses with <code>rwuniq --fields</code> . . . . .	66
4.9	Profiling IP addresses with <code>rwstats --fields</code> . . . . .	67
4.10	Isolating DNS and Non-DNS Behavior with <code>rwuniq</code> . . . . .	68
4.11	Generating Bags with <code>rwbag</code> . . . . .	68
4.12	Summarizing Network Traffic with <code>rwuniq</code> . . . . .	69
4.13	Summarizing Network Traffic with Bags . . . . .	69
4.14	Creating a Bag of Network Scanners with <code>rwbagbuild</code> and <code>rwscan</code> . . . . .	71
4.15	Viewing the Contents of a Bag with <code>rwbagcat</code> . . . . .	71
4.16	Thresholding Results with <code>rwbagcat --mincounter</code> , <code>--maxcounter</code> , <code>--minkey</code> , and <code>--maxkey</code> . . . . .	72
4.17	Displaying Unique IP Addresses per Value with <code>rwbagcat --bin-ips</code> . . . . .	72
4.18	Displaying Decimal and Hexadecimal Output with <code>rwbagcat --key-format</code> . . . . .	73
4.19	Creating an IP Set from a Bag with <code>rwbagtool --coverset</code> . . . . .	74
4.20	Using <code>rwbagtool --intersect</code> to Extract a Subnet . . . . .	75
4.21	Abstracting Source IPv4 addresses with <code>rwnetmask</code> . . . . .	75

4.22	Generating a Monitored Address Space IPset with <code>rwsetbuild</code>	76
4.23	Generating a Broadcast Address Space IPset with <code>rwsetbuild</code>	76
4.24	Performing an IPset Union with <code>rwsettool</code>	77
4.25	Displaying Repository Dates with <code>rwsiteinfo</code>	77
4.26	Counting Outbound DNS Servers with <code>rwset</code>	77
4.27	Finding IPset Differences with <code>rwsettool</code>	78
4.28	Finding IPset Symmetric Difference with <code>rwsettool</code>	78
4.29	Grouping Outbound DNS Servers by Sensor	79
4.30	Identifying DNS Traffic Flow	79
4.31	Identifying Shared DNS Monitoring	80
4.32	Displaying the Contents of IP Sets with <code>rwsetcat</code>	80
4.33	<code>rwsetcat</code> Options for Showing Structure	81
4.34	Grouping Flows of a Long Session with <code>rwgroup</code>	84
4.35	Dropping Trivial Groups with <code>rwgroup --rec-threshold</code>	85
4.36	Summarizing Groups with <code>rwgroup --summarize</code>	86
4.37	Using <code>rwgroup</code> to Identify Specific Sessions	87
4.38	Using <code>rwmatch</code> with Incomplete Relate Values	88
4.39	Using <code>rwmatch</code> with Full TCP Fields	89
4.40	<code>rwfileinfo</code> for Sets, Bags, and Prefix Maps	91
5.1	Building an IPset Inventory for Sensor S0	94
5.2	Automating IPset Inventories	98
6.1	Using <code>rwfilter</code> to Profile NTP Activity	102
6.2	Using <code>rwuniq</code> to examine NTP Activity	103
6.3	Using <code>rwcount</code> to generate NTP Timelines	103
6.4	Using <code>rwuniq</code> and Bags to Summarize Prior Traffic on NTP Clients	105
6.5	Using Multiple Data Pulls to Filter on Multiple Criteria	107
6.6	Filtering on Multiple Criteria with a Tuple File	109
6.7	Merging the Contents of Bags Using <code>rwbagtool --add</code>	111
6.8	Using <code>rwbagtool</code> to Generate Percentages	113
6.9	Using <code>rwset</code> to Filter for a Set of Scanners	113
6.10	Using <code>rwbagtool</code> to Filter Out a Set of Scanners	115
6.11	Combining Flow Record Files with <code>rwcat</code> to Count Overall Volumes	117
6.12	<code>rwsplit</code> for Coarse Parallel Execution	118
6.13	<code>rwsplit</code> to Generate Statistics on Flow Record Files	119
6.14	Simple File Anonymization with <code>rwtuc</code>	120
6.15	Using <code>rwmapbuild</code> to Create a FCC Pmap File	123
6.16	Using Pmap Parameters with <code>rwfilter</code>	126
6.17	Viewing Prefix Map Labels with <code>rwcut</code>	127
6.18	Sorting by Prefix Map Labels	127
6.19	Counting Records by Prefix Map Labels	128
6.20	Query Addresses and Protocol/Ports with <code>rwmaplookup</code>	129
7.1	Looking for Service Ports with Higher Inbound than Outbound TCP Traffic	134
7.2	Identifying Abnormal TCP Flows and their Originating Hosts	136
7.3	Finding Activity of Illegitimate Destination IP Addresses	138
7.4	Finding Changed Behavior in Destination IPs	139
8.1	<code>ThreeOrMore.py</code> : Using PySiLK for Memory in <code>rwfilter</code> Partitioning	144
8.2	<code>portknock.py</code> : Using PySiLK to Retain State in <code>rwfilter</code> Partitioning	146
8.3	Calling <code>ThreeOrMore.py</code>	147
8.4	Using <code>--python-expr</code> for Partitioning	148

8.5	<code>vpn.py</code> : Using PySiLK with <code>rwfilter</code> for Partitioning Alternatives	148
8.6	<code>matchblock.py</code> : Using PySiLK with <code>rwfilter</code> for Structured Conditions	150
8.7	Calling <code>matchblock.py</code>	151
8.8	<code>delta.py</code>	152
8.9	Calling <code>delta.py</code>	153
8.10	<code>payload.py</code> : Using PySiLK for Conditional Fields with <code>rwsort</code> and <code>rwcut</code>	154
8.11	Calling <code>payload.py</code>	155
8.12	<code>decode_duration.py</code> : A Program to Create a String Field for <code>rwcut</code>	155
8.13	Calling <code>decode_duration.py</code>	156
8.14	<code>sitefield.py</code> : A Program to Create a String Field for Five SiLK Tools	157
8.15	Calling <code>sitefield.py</code>	158
8.16	<code>bpp.py</code>	158
8.17	Calling <code>bpp.py</code>	159
B.1	A UNIX Command Prompt	173
B.2	Using Simple UNIX Commands	175
B.3	Output Redirection	176
B.4	Input Redirection	177
B.5	Using a Pipe	177
B.6	Using a Here-Document	178
B.7	Using a Named Pipe	179



# Acknowledgements

The authors wish to acknowledge the valuable contributions of all members of the CERT® Situational Awareness group and the CERT Engineering Group, past and present, to the concept and execution of the SiLK tool suite and to this handbook. Many individuals served as contributors, reviewers, and evaluators of the material in this handbook.

The authors also gratefully acknowledge the many SiLK users who have contributed immensely to the evolution of the tool suite.

Lastly, the authors wish to acknowledge their ongoing debt to the memory of Suresh L. Konda, PhD, who led the initial concept and development of the SiLK tool suite as a means of gaining network situational awareness.



# Handbook Goals

## How to Use This Handbook

This handbook is an introduction to methods of analyzing network traffic, illustrated by commands from the SiLK tool suite. The focus is on learning to identify traffic features important to the security of information on the network. The handbook moves from a basic understanding of network flow and the SiLK tool suite through a series of examples that illustrate how to use SiLK to analyze network behavior.

The examples in this handbook are mainly command sequences that illustrate specific analysis concepts. Examples are commonly discussed on a line-by-line basis in the text and presented as command and output listings. In general, examples are also associated with a specific task (or tasks), indicated in the section and in the example caption. Case studies take a deeper dive into specific topics for analysis.

For readers already familiar with SiLK, the explanations of SiLK commands in the text of this handbook are kept short enough not to be redundant. More complete discussion of the commands and their parameters are provided in the appendices of this guide, the *SiLK Reference Guide*, and the `man` pages for the SiLK commands. Readers who are interested in analyzing network flow records with other tools than SiLK are encouraged to read the overall description of the analysis approaches, then use the description of commands to find parallels using the tool suite of their choice.

## How This Handbook Is Organized

This handbook contains the following chapters:

1. **Introduction to SiLK** provides a short overview of some of the background necessary to begin using the SiLK tools for analysis. It includes a brief introduction to the SiLK suite and describes the basics of network flow capture by sensors and storage in the SiLK flow repository. It also discusses the analysis process used in this handbook.
2. **Basic Single-path Analysis with SiLK: Profiling and Reacting** describes the the most straightforward analysis approach and applies it to several example analyses. It introduces some of the core SiLK commands and uses them to analyze network traffic.
3. **Case Studies: Basic Single-path Analysis** applies the single-path analysis approach to several extended examples, focusing on how those examples were developed from an initial problem statement through executable commands.
4. **Intermediate Multi-path Analysis with SiLK: Explaining and Investigating** explains a more complex, intermediate form of analysis which applies basic, single-path analysis in a multi-pronged

structure. The chapter describes how multi-path analysis can be applied and includes a fuller exploration of SiLK tools that may be useful for this type of analysis.

5. **Case Studies: Intermediate Multi-path Analysis** applies multi-path analysis to extended examples.
6. **Advanced Exploratory Analysis with SiLK: Exploring and Hunting** discusses the use of SiLK to deal with open-ended, often iterative analyses that incorporate both single-path and multi-path methods. It also describes more sophisticated uses of the SiLK tool suite that support complex analyses of network behavior.
7. **Case Studies: Advanced Exploratory Analysis** applies exploratory analysis to an extended example.
8. **Extending the Reach of SiLK with PySiLK** describes how to extend the functionality of the SiLK tool suite by using the Python scripting language.

The appendices to this guide introduce fundamental networking concepts, describe useful Unix commands, summarize the SiLK commands referenced in this guide, and list sources for additional information about the SiLK tool suite and network analysis.

## What This Handbook Doesn't Cover

This handbook does not contain an exhaustive description of all the tools in the SiLK tool suite or of all the options in the described commands. Rather, it offers concepts and examples to allow analysts to accomplish needed work while continuing to build their skills and familiarity with SiLK.

- Every SiLK tool includes a `--help` option that briefly describes the command and lists its parameters.
- Every tool also has a manual page (also called a `man` page) that provides detailed information about the use of the tool. These pages may be available on your system by typing `man command`; for example, `man rfilter` to see information about the `rfilter` command.
- The SiLK Documentation page at <https://tools.netsa.cert.org/silk/docs.html> includes links to individual manual pages.
- The *SiLK Reference Guide* is a single document that bundles all of the SiLK manual pages. It is available in HTML and PDF formats on the SiLK Documentation page (<https://tools.netsa.cert.org/silk/docs.html>).
- Various analysis topics are explored via tooltips, available at <https://tools.netsa.cert.org/tooltips.html>.

This handbook deals solely with the analysis of network flow record data using an existing installation of the SiLK tool suite. For information on installing and configuring a new SiLK tool setup and on the collection of network flow records for use in these analyses, see the “Installation Information” section of the SiLK Documentation page at <https://tools.netsa.cert.org/silk/docs.html#installation>.



# Chapter 1

## Introduction to SiLK

Network analysts need to build an ongoing perspective on the traffic passing over their networks. This perspective is often built on information about the traffic (such as volumes, timing, and communication paths), rather than on the traffic itself. This chapter introduces the tools and techniques used to store such information, particularly in the form known as *network flow*. It will help you to become familiar with the structure of network flow data, how the SiLK collection system gathers those data from sensors, and how to use those data.

Upon completion of this chapter you will be able to

- describe a network flow record and the conditions under which the collection of one begins and ends
- describe the types of SiLK flow records
- describe the structure of the SiLK flow repository
- understand the steps involved in analyzing network flow data
- describe the dataset for the examples in this guide

### 1.1 What is SiLK?

The System for internet-Level Knowledge<sup>1</sup> (SiLK) tool suite is a highly scalable flow-data capture and analysis system developed by the CERT Situational Awareness group at Carnegie Mellon University’s Software Engineering Institute (SEI). The SiLK tools provide network security analysts with the means to understand, query, and summarize both recent and historical traffic data represented as network flow records (also referred to as “network flow” or “network flow data” and occasionally just “flow”). These tools provide network security analysts with a relatively complete high-level view of traffic across an enterprise network, subject to placement of sensors.

Analyses using the SiLK tools provide insight into various aspects of network behavior. Some example applications of this tool suite include:

---

<sup>1</sup>The suite name, and in particular the capitalization, were chosen in memory of Dr. Suresh L. Konda, who was the inspirational leader for the creation of the initial suite prior to his sudden passing.

- supporting network forensics: identifying artifacts of intrusions, vulnerability exploits, worm behavior, etc.
- providing service inventories for large and dynamic networks (on the order of a /8 Classless Inter-Domain Routing (CIDR) block)
- generating profiles of network usage (bandwidth consumption) based on protocols and common communication patterns
- enabling non-signature-based scan detection and worm detection, for detection of limited-release malicious software and for identification of precursors

These examples, and others, are explained further in this handbook. By providing a common basis for these analyses, the SiLK tools provide a framework for developing network situational awareness.

Common questions addressed via flow analyses include (but aren't limited to)

- What is on my network?
- What constitutes typical network behavior?
- What happened before, during, and after an event?
- Where are policy violations occurring?
- Which are the most popular web servers?
- How much volume would be reduced by applying a blacklist?
- Do my users browse to known infected web servers?
- Is a spammer on my network?
- When did my web server stop responding to queries?
- Is my organization routing undesired traffic?
- Who uses my public Domain Name System (DNS) server?

## 1.2 The SiLK Flow Repository

### 1.2.1 What is Network Flow Data?

NetFlow is a traffic-summarizing format that was first implemented by Cisco Systems® primarily for accounting purposes. Network flow data (or network flow) is a generalization of NetFlow. Network flow collection differs from direct packet capture (such as with `tcpdump`) in that it builds a summary of communications between sources and destinations on a network. For NetFlow, this summary covers all traffic matching seven relevant keys: the source and destination IP addresses, the source and destination ports, the transport layer protocol, the type of service, and the router interface.

SiLK uses five of these attributes to constitute the *flow label*:

1. source IP address

2. destination IP address
3. source port
4. destination port
5. transport layer protocol

These attributes (also known as the *five-tuple*), together with the start time of each network flow, distinguish network flows from each other. The SiLK *repository* stores the accumulated flows from a network.

### 1.2.2 Structure of a Flow Record

A network flow often covers multiple packets that all match the fields of their common labels. A *flow record* thus provides the label and statistics on the packets covered by the network flow, including the number of packets covered by the flow, the total number of bytes, and the duration and timing of those packets (among other fields). A *flow file* is a series of flow records.

The fields in the flow record are listed in Table 1.1. Every field is identified by a name and number that can be used interchangeably. For example, the source IP address field of a flow record can be identified by either its field name (`sIP`) or its field number (1). Capitalization does not matter: `sIP` is equivalent to `sip` or `SIP`.

Because network flow is a summary of traffic, it does not contain packet payload data, which are expensive to retain on a large, busy network. Each network flow record created by SiLK is very small: it can be as little as 22 bytes (the exact size is determined by several configuration parameters). However, even at that tiny size, a sensor may collect many gigabytes of flow records daily on a busy network.

Some of the fields are actually stored in the record, such as start time and duration. Some fields are not actually stored; rather, they are derived either wholly from information in the stored fields or from a combination of fields stored in the record and external data. For example, end time is derived by adding the start time and the duration. Source country code is derived from the source IP address and a table that maps IP addresses to country codes.

### 1.2.3 Flow Generation and Collection

To understand how to use SiLK for analysis, it helps to have some knowledge of how network flow data are collected, stored, and managed. Understanding how the data are partitioned can produce faster queries by reducing the amount of data searched. In addition, by understanding how the sensors complement each other, it is possible to gather traffic data even when a specific sensor has failed.

Every day, SiLK may collect many gigabytes of network flow records from across the enterprise network. This section reviews the collection process and shows how data are stored as network flow records.

A network flow record is generated by sensors throughout the enterprise network. Usually, the majority of these sensors are routers. Specialized sensors such as `yaf`<sup>2</sup> can be employed when a data feed from a router is not available, such as on a home network or on an individual host. `yaf` can also be used to avoid artifacts in a router's implementation of network flow or to use non-device-specific network flow data formats such

---

<sup>2</sup><https://tools.netsa.cert.org/yaf/>

Table 1.1: Fields in a SiLK Network Flow record

Field Number	Field Name	Description
1	sIP	Source IP address for flow
2	dIP	Destination IP address for flow
3	sPort	Source port for flow (or 0)
4	dPort	Destination port for flow (or 0)
5	protocol	Transport layer protocol number for flow
6	packets, pkts	Number of packets in flow
7	bytes	Number of bytes in flow (starting with IP header)
8	flags	Cumulative TCP flag fields of flow (or blank)
9	sTime	Start date and time of flow
10	duration	Duration of flow
11	eTime	End date and time of flow
12	sensor	Sensor that collected the flow
13	in	Ingress interface or VLAN on sensor (usually zero)
14	out	Egress interface or VLAN on sensor (usually zero)
15	nhIP	Next-hop IP address (usually zero)
16	sType	Type of source IP address (pmap required)
17	dType	Type of destination IP address (pmap required)
18	scc	Source country code (pmap required)
19	dcc	Destination country code (pmap required)
20	class	Class of sensor that collected flow
21	type	Type of flow for this sensor class
—	iType	ICMP type for ICMP and ICMPv6 flows (SiLK V3.8.1+)
—	iCode	ICMP code for ICMP and ICMPv6 flows (SiLK V3.8.1+)
25	icmpTypeCode	Both ICMP type and code values (before SiLK V3.8.1)
26	initialFlags	TCP flags in initial packet
27	sessionFlags	TCP flags in remaining packets
28	attributes	Termination conditions
29	application	Standard port for application that produced the flow

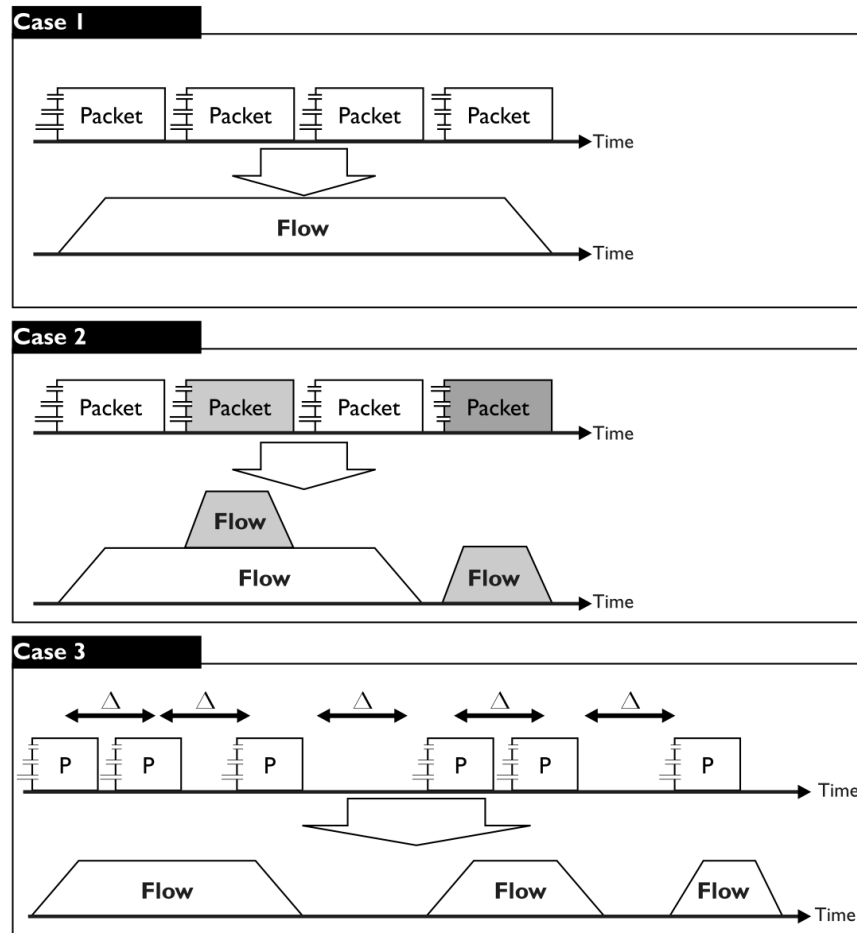
as IPFIX<sup>3</sup>. It provides more control over network flow record generation and can convert packet data to network flow records via a script that automates this process.

A sensor generates network flow records by grouping together packets that are closely related in time and have a common flow label. “Closely related” is defined by the sensor and typically set to around 30 seconds. Figure 1.1 shows the generation of flows from packets. Case 1 in that figure diagrams flow record generation when all the packets for a flow are contiguous and uninterrupted. Case 2 diagrams flow record generation when several flows are collected in parallel. Case 3 diagrams flow record generation when timeout occurs, as discussed below.

Network flow is an approximation of traffic. Routers and other sensors make a guess when they decide which packets belong to a flow. These guesses are not perfect; there are several well-known phenomena in which a

<sup>3</sup>See <https://tools.ietf.org/html/rfc7011> for definitions of the IPFIX information elements; see the IPFIX protocol description and <https://www.iana.org/assignments/ipfix> for their descriptions.

Figure 1.1: From Packets to Flows



long-lived session will be split into multiple flow records:

1. *Active timeout* is the most common cause of a split network flow. Network flow records are purged from the sensor's memory and restarted after a configurable period of activity. As a result, all network flow records have an upper limit on their duration that depends on the local configuration. A typical value would be around 30 minutes.
2. *Cache flush* is a common cause of split network flows for router-collected network flow records. Network flows take up memory resources in the router, and the router regularly purges this cache of network flows for housekeeping purposes. The cache flush takes place approximately every 30 minutes as well. A plot of network flows over a long period of time shows many network flows terminate at regular 30-minute intervals, which is a result of the cache flush.
3. *Router exhaustion* also causes split network flows for router-collected flows. A router has limited processing and memory resources devoted to network flow. During periods of stress, the flow cache will fill and empty more frequently due to the number of network flows collected by the router.

Use of specialized flow sensors can avoid or minimize cache-flush and router-exhaustion issues. All of these cases involve network flows that are long enough to be split. As we show later, the majority of network flows collected at the enterprise network border are small and short-lived.

## 1.2.4 Introduction to Flow Collection

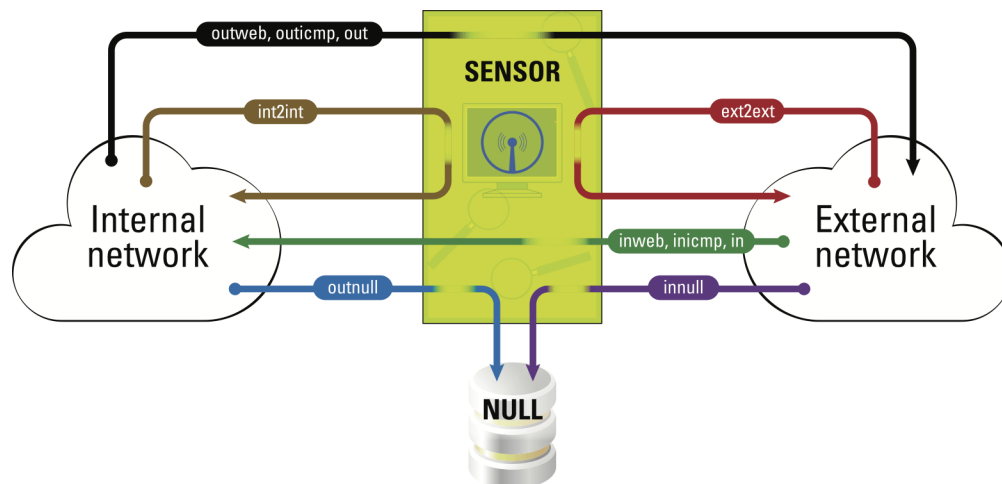
An enterprise network comprises a variety of organizations and systems. The flow data to be handled by SiLK are first processed by the collection system, which receives flow records from the sensors and organizes them for later analysis. The collection system may collect data through a set of sensors that includes both routers and specialized sensors that are positioned throughout the enterprise network. After records are added to the flow repository by the collection system, analysis is performed using a custom set of software called the SiLK analysis tool suite.

The SiLK project is active, meaning that the system is continually improved. These improvements include new tools and revisions to existing collection and analysis software. See Appendix E for information on how to obtain the most up-to-date version of SiLK.

## 1.2.5 Where Network Flow Data Are Collected

While complex networks may segregate flow records based on where the records were collected (e.g., the network border, major points within the border, at other points), the generic implementation of the SiLK collection system defaults to collection only at the network border, as shown in Figure 1.2. The default implementation has only one class of sensors: *all*. Further segregation of the data is done by type of traffic.

Figure 1.2: Default Traffic Types for Sensors



The SiLK tool `rwsiteinfo` can produce a list of sensors in use for a specific installation, reflecting its configuration. For more information on how to use this tool, see Section 2.2.1.

### 1.2.6 Types of Enterprise Network Traffic

In SiLK, the term *type* mostly refers to the direction of traffic, rather than a content-based characteristic. In the generic implementation (as shown in Figure 1.2), there are six basic types and five additional types. The basic types are

- **in** and **inweb**, which is traffic coming from the Internet service provider (ISP) to the enterprise network through the border router. Web traffic is separated from other traffic due to its volume, making many searches faster.
- **out** and **outweb**, which is traffic coming from the enterprise network to the ISP through the border router.
- **int2int**, which is traffic going both from and to the enterprise network, but which passes by the sensor.
- **ext2ext**, which is traffic going both from and to the ISP, but which passes by the sensor. (The presence of this type of traffic usually indicates a configuration problem either in the sensor or at the ISP.)

The additional SiLK types are

- **inicmp** and **outicmp**, which represent ICMP traffic entering or leaving the enterprise network. These types are operational only if SiLK was compiled with the option to support them.
- **innull** and **outnull**, which only can be found when the sensor is a router and not a dedicated sensor. They represent traffic from the upstream ISP or the enterprise network, respectively, that terminates at the router's IP address or is dropped by the router due to an access control list.
- **other**, which is assigned to traffic where one of the addresses (source or destination) is in neither the internal nor the external networks.
- The constructed type **all** selects all types of flows associated with a class of sensors.

These types are configurable. Configurations vary as to which types are in actual use (see the discussion below under [Sensors: Class and Type](#)).

### 1.2.7 The Collection System and Data Management

Data collection starts when a flow record is generated by one of the sensors: either a router or a dedicated sensor. Flow records are generated when a packet relevant to the flow is seen, but a flow is not *reported* until it is complete or flushed from the cache. Consequently, a flow can be seen some time after the start time of the first packet in the flow, depending on timeout configuration and on sensor caching, among other factors.

Packed flows are stored into files indicated by class, type, sensor, and the hour in which the flow started. So for traffic coming from the ISP through or past the sensor named SEN1 on March 1, 2018 for flows starting between 3:00 and 3:59:59.999 p.m. Coordinated Universal Time (UTC), a sample path to the file could be `/data/SEN1/in/2018/03/01/in-SEN1_20180301.15`.

## 1.2.8 How Network Flow Data Are Organized

The data repository is accessed via the SiLK tools, particularly the `rwfilter` command. An analyst uses `rwfilter` to choose the type of data to be viewed by specifying a set of selection parameters. This handbook discusses selection parameters in more detail in Section 2.2.2 and Appendix C.3; this section briefly outlines how data are stored in the repository.

### Dates

The SiLK repository stores data in hourly divisions, which are referred to in the form `yyyy/mm/ddThh` in UTC. Thus, the hour beginning 11 a.m. on February 23, 2018 in Pittsburgh would be referred to as `2018/2/23T16` when compensating for the difference between UTC and Eastern Standard Time (EST)—five hours.

In general, data for a particular hour starts being recorded at that hour and will continue recording until some time after the end of the hour. Under ideal conditions, the last long-lived flows will be written to the file soon after they time out (e.g., if the active timeout period is 30 minutes, the last flows will be written out 30 minutes plus propagation time after the end of the hour). Under adverse network conditions, however, flows could accumulate on the sensor until they can be delivered. Under normal conditions, the file for `2018/3/7 20:00 UTC` would have data starting at 3 p.m. in Pittsburgh and finish being updated after 4:30 p.m. in Pittsburgh.

### Sensors: Class and Type

Data are divided by time and sensor. The class of a sensor is often associated with the sensor's role as a router: access layer, distribution layer, core (backbone) layer, or border (edge) router. The classes of sensors that are available are determined by the installation. By default, there is only one class—`all`—but based on analytical interest, other classes may be configured as needed. As shown in Figure 1.2, each class of sensor has several types of traffic associated with it: typically `in`, `inweb`, `out`, and `outweb`.

Data types are used for two reasons:

1. They group data together into common directions.
2. They split off major query classes.

As shown in Figure 1.2, most data types have a companion web type (i.e., `in` and `inweb`, `out` and `outweb`). Web traffic generally constitutes about 50% of the flows in any direction; by splitting the web traffic into a separate type, we reduce query time.

Most queries to repository data access one *class* of data at a time but access multiple *types* simultaneously.

## 1.3 The SiLK Tool Suite

The SiLK analysis suite consists of over 60 command-line UNIX tools (including flow collection tools) that rapidly process flow records or manipulate ancillary data. The tools can communicate with each other and with scripting tools via pipes (both unnamed and named) or via intermediate files; see Section B.2 for more information.



Flow analysis is generally input/output bound—the amount of time required to perform an analysis is proportional to the amount of data read from disk. A major goal of the SiLK tool suite is to minimize that access time. Some SiLK tools perform functions analogous to common UNIX command-line tools and to higher level scripting languages such as Perl.<sup>®</sup> However, the SiLK tools process this data in non-text (binary) form and use data structures specifically optimized for analysis.

Consequently, most SiLK analysis consists of a sequence of operations using the SiLK tools. These operations typically start with an initial `rwfilter` call to retrieve data of interest and culminate in a final call to a text output tool like `rwstats` or `rwuniq` to summarize the data for presentation.

Keeping data in binary for as many steps as possible greatly improves efficiency of processing. This is because the structured binary records created by the SiLK tools are readily decomposed without parsing, their fields are compact, and the fields are already in a format that is ready for calculations, such as computing netmasks.

In some ways, it is appropriate to think of SiLK as an awareness toolkit. The flow-record repository provides large volumes of data, and the tool suite provides the capabilities needed to process these data. However, the actual insights come from analysts.

## 1.4 How to Use SiLK for Analysis

The SiLK tool suite provides a robust collection of tools to facilitate network traffic analysis tasks. It is designed to be very flexible in its support of analysis methods. Over time, different analysts have used a variety of approaches in their use of SiLK. This section discusses three approaches that have been useful in analyzing network flow records.

The chapters following this one expand on these approaches in more detail, focusing on the support that network flow analysis can provide to such analyses. Being aware of and practicing multiple approaches to analysis enables an analyst to gain insight into a wide variety of network traffic behaviors.

### 1.4.1 Single-path Analysis

The *single-path* approach is the most basic and most commonly-used approach to analyzing network behavior. It makes use of a single sequence of commands to produce the analytic results. In this approach, the analyst formulates an initial hypothesis, constructs a query to retrieve traffic of interest, produces a table, summary, or series to profile this traffic, and then interprets this profile either numerically or through a graph. Iteration can be used if needed (e.g., to refine the initial query), but may not be necessary for many simpler, more straightforward analyses.

This approach could be used for service identification, network device inventories, incident response, or usage studies. Chapter 2 provides an overview of single-path analysis, including the SiLK commands that are most commonly used with it. Chapter 3 describes example case studies of single-path analyses.

### 1.4.2 Multi-path Analysis

The *multi-path* approach uses a sequence of tools that frequently involve several alternatives, and often includes iterating over some steps. Although a multi-path approach can be done manually, it more often involves scripting to select alternatives based on categories of data and then iterate until the desired traffic is isolated or the desired summaries are produced. The alternatives are used as required for processing groups of records in differing ways to reach results that profile behavior of interest.

This approach could be used for examining traffic using several protocols, each following its own alternative set of characteristics, to accomplish the same goal. For example, there are multiple ways that malware can beacon to its command-and-control network. Each of those ways could be examined separately via a chain of SiLK commands, generating sets of results that contribute to overall awareness of beaconing.

Chapter 4 provides an overview of multi-path analysis, including the SiLK commands that are most commonly used with it. Chapter 5 describes an example case study of multi-path analysis.

### 1.4.3 Exploratory Analysis

We do not always know ahead of time what the scope of our analysis will be—or even what questions we should be asking! *Exploratory* analysis is an open-ended approach to formulating, scoping, and conducting a network analysis. It uses single-path and multi-path analyses as building blocks for investigating anomalous network traffic. These simpler types of analysis help us to formulate different scenarios, investigate alternative hypotheses, and explore multiple aspects of the data. Exploratory analysis is initially manual in nature, but can transition to scripted analysis for ease of repetition and for regularity of results.

This approach is used for complex or emerging phenomena, where multiple indicators need to be combined to gain understanding. An example of this approach to analysis would be a study of data exfiltration, which can be performed in a wide variety of ways. Each of those exfiltration methods could be profiled using a set of indicators, and the results of all such analyses combined to produce a composite understanding of traffic being passed to various groups of suspicious addresses.

Chapter 6 provides an overview of exploratory analysis, including advanced SiLK commands and concepts. Chapter 7 describes an example case study of exploratory analysis.

## 1.5 Workflow for SiLK Analysis

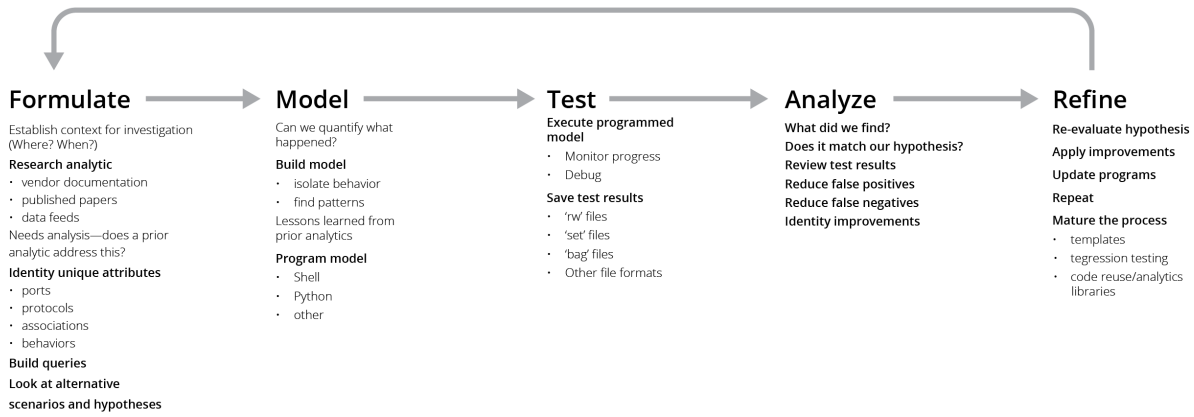
SiLK analyses share a common workflow, shown in Figure 1.3. While single-path, multi-path, and exploratory analysis may incorporate different steps in this workflow, all follow its general sequence.

### 1.5.1 Formulate

The *Formulate* step investigates the context of the event. Essentially, it involves collecting information to identify the unique attributes of the network, its operation, and the event. How large is the network? How is it structured? Where are network sensors located? When did the event occur? Is it associated with specific sensors, IP addresses, hosts, network spaces, ports, protocols, and so forth? Do any earlier analyses of the network offer insight? The information may be incomplete at this point, but it serves as a starting point for launching the analysis and establishing its scope. We can use it to formulate a hypothesis for the network's behavior. This hypothesis serves as the basis of our analysis. In more sophisticated exploratory analyses, we can formulate multiple scenarios and hypotheses for investigation and analysis.

Information gleaned from exploring the event's context helps us to establish which network behaviors should be included in (or excluded from) our analysis. We can use this information to construct a *query* to select and partition network flow records from the SiLK repository or a stored file. Queries typically incorporate information such as where the flow was collected, the date of data collection, and the flow direction. Within the SiLK community, query selection is commonly called a *data pull*.

Figure 1.3: SiLK Analysis Workflow



*Partitioning* applies tests to selected flow records to separate (or partition) them into categories for further inspection and investigation. A default set of tests is provided with SiLK. It includes IP addresses, ports, protocols, time, and volumes. (If additional tests are needed for analyses, the SiLK tools can be extended via *plugins* to provide them.)

The combination of selection and partitioning (commonly referred to as *filtering*) is performed with the `rwfilter` command. Records that meet the filtering criteria are sent to *pass* destinations. Records that do not are sent to *fail* destinations. Both can be combined into *all* destinations. This provides flexible options to either store query results in files or use pipes to send them to other commands for processing.

## 1.5.2 Model

The *Model* step summarizes data and investigates behaviors of interest. What is the network's behavior during normal operation? What happened during an event? What patterns and behaviors can we identify? Are they similar to those observed during other events? By examining the information gathered during the Formulate step, you can come up with a model of the event that perhaps explains what is going on.

SiLK provides a variety of tools for examining network flow data associated with an event. Each tool offers different views into the data that can be considered independently or in combination for analysis. For example, SiLK includes tools for generating time-series summaries of traffic (the `rwcount` command), computing summary statistics (the `rwstats` command), and summing up the values of flow attributes for user-defined time periods (the `rwuniq` command).

This step can be done manually. For analyses that are larger in scope, it can be automated by using shell or Python scripts.

### 1.5.3 Test

The *Test* step runs the model that you created—either manually or by executing shell or Python scripts. This gives you a chance to check the progress of the analysis.

SiLK includes commands for sorting flow records according to user-defined keys (the `rwsort` command), creating sets of unique IP addresses from flow records (`rwset` and its related commands), and creating groups of records by other criteria (`rwbag` and its related commands). These commands help you to organize output from the various SiLK commands and save it for further use.

### 1.5.4 Analyze

The *Analyze* step reviews the results of the previous steps. What do these results tell us about the event? What behaviors have been identified? What types of events are they associated with? What relationships can we identify between flows? Do our initial hypotheses still hold up? Can we find and eliminate false positives and false negatives?

This step involves examining and interpreting output from the analysis tools mentioned earlier. SiLK can also translate binary flow records into text for analysis with graphics packages, spreadsheets, and mathematical tools (the `rwcut` command).

### 1.5.5 Refine

The *Refine* step improves the analysis. Did we successfully explain the event? If not, what problems did we encounter? Did we properly understand the event's context? Did our query into the SiLK repository pull too much data? Do we need to dig deeper into the data during the modeling and testing steps? Should we take another look at the results to see if we missed or misinterpreted important patterns and behaviors?

The preceding steps in the workflow can be combined in an iterative pattern. For example, you may want to isolate flow records of interest from unrelated network traffic by making additional queries with the `rwfilter` command and repeating subsequent steps in the analysis. This narrows the data to focus on the time periods and behaviors of interest and eliminate unneeded flow records.

The workflow described in this section gives us the flexibility to begin our data exploration with a general question, apply one or more analyses to the question, and complete the workflow with a repeatable analytic. This flexibility does come with trade-offs, however. Queries typically increase proportionally with the time window and flow record attributes of an analysis. Therefore, a precise model of an analysis should be produced to minimize the query results.

## 1.6 Applying the SiLK Workflow

The SiLK workflow can be applied in different ways to meet the requirements of analysis groups. Groups that are primarily concerned with network operations will often focus on network monitoring or service and device validation. Incident response groups commonly focus on changes in network behavior that may be associated with an incident. Security improvement groups often focus on understanding problematic network behavior and changes that identify the impact of improvements.

While the SiLK suite offers features that support all groups, the work required to use them will vary. Many of these applications are addressed in the remaining chapters of this handbook.

## 1.7 Dataset for Single-path, Multi-path, and Exploratory Analysis Examples

The dataset used for the command examples and case studies of single-path, multi-path, and exploratory analysis in this document is the FCCX-15 dataset<sup>4</sup>. It originates from a June, 2015 Cyber Exercise conducted by the Software Engineering Institute at Carnegie Mellon University in a virtual environment.

The exercise network topology is documented in the data download and comprises a distributed enterprise for the period from June 2-16, 2015. Internet and transport layer protocols such as IPv4, TCP, UDP, and ICMP are well represented in the data. Link layer protocols such as IGMP and OSPF are also included; however, they are not as prevalent as the Internet and transport layer protocols.

The analyses in this guide investigate events that occurred during this exercise.

---

<sup>4</sup><https://tools.netsa.cert.org/silk/referencedata.html>



## Chapter 2

# Basic Single-path Analysis with SiLK: Profiling and Reacting

This chapter introduces basic single-path analysis through application of the analytic development process with the SiLK tool suite. It discusses basic analysis of network flow data with SiLK, in addition to specific tools and how they can be combined to form a workflow.

Upon completion of this chapter you will be able to

- describe basic single-path analysis and how it maps to the analytic development process
- understand SiLK tools commonly used with basic single-path analysis
- describe SiLK IPsets and their application
- describe the single-path analysis workflow using network flow data

### 2.1 Single-path Analysis: Concepts

Single-path analysis is the approach of combining data with methods that do not require conditional steps, integration, or a great deal of refinement. In layman’s terms, single-path analysis can be described as the ‘start-to-finish’ approach of combining one or more analytical steps to characterize network behavior. Its output may contain multiple attributes and characteristics; however, it results in information that normally does not need continued iteration. Figure 2.1 provides an overview of single-path analysis.

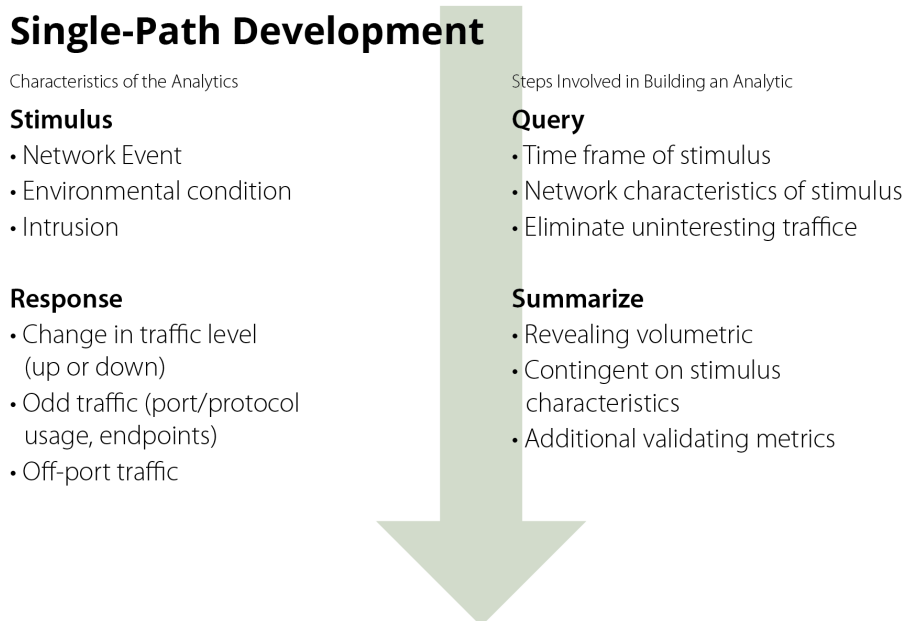
Single-path analysis typically incorporates the Formulate, Model, Test, and Analyze steps of the analysis workflow described in Section 1.5. The Refine step can also be included—for instance, to change the scope of a data pull from the SiLK repository—but is not always needed. It begins by identifying the context of an analysis and formulating a hypothesis to explain the behavior under investigation. Event attributes such as hosts, networks, and time periods are used to identify, retrieve, and partition data for analysis. Attributes such as frequency, volume, and supporting network services provide additional behavioral context.

Single-path analysis then summarizes this data to produce sequences of event behavior. Data can be separated into logical groups such as successful and unsuccessful contacts, scanning, and misconfiguration. This

enables analysts to produce known and unknown activity, trends, and differences for comprehensive analysis of a network's behavior.

Analysts also use single-path analyses as building blocks for broader, more complex analyses. See Chapter 4 and Chapter 6 for descriptions of analysis workflows that include single-path analyses as building blocks for more comprehensive investigations of network activity.

Figure 2.1: Single-Path Analysis



### 2.1.1 Scoping Queries of Network Flow Data

Be careful when defining the scope of an initial query into the SiLK repository or a file. The natural tendency is to make the partitioning criteria very inclusive, which has two drawbacks. It pulls over-large amounts of data, consuming storage and other computer resources. Overly-broad queries may match behaviors other than those of interest, which will complicate later steps in the analysis.

The preferred method for scoping queries is the opposite:

1. Make the partitioning criteria initially narrow, specific to the desired behaviors.
2. Once the traffic related to the behavior is retrieved, broaden the initial criteria to identify related network traffic.

Starting narrow and broadening the scope of the data query as the analysis proceeds will use computing resources more efficiently and facilitate clearer analysis by minimizing unwanted data.



### 2.1.2 Excluding Unwanted Network Traffic

Despite narrowing the initial query, unrelated traffic is commonly included in the initially-retrieved records. Analysis will often require isolating the desired activity from among the retrieved traffic. This may involve studying the traffic to identify unrelated behaviors and constructing further criteria to exclude them. It may also include eliminating traffic involving specific addresses (often by using the `rwset` tools to build IPsets), traffic that does not occur in the proper timeframe (often by using a further `rwfilter` call), or traffic that lacks specific protocol information associated with behavior of interest (again by using a further `rwfilter` call).

### 2.1.3 Example Single-Path Analysis

This chapter documents an example single-path analysis using the SiLK tool suite. It serves as an independent analysis, but could also depict the beginning of a multi-path (intermediate) or exploratory (advanced) analysis workflow. The example begins with identifying the context of an event by using `rwsiteinfo` to select relevant sensors and time periods for analysis. The time window is expanded beyond the event under analysis to select data to compare against the event period. `rwfilter` is then used to retrieve network flow records that apply to the defined sensor and period.

Traffic characteristics such as bytes, packets, and TCP flags options are then used with `rwfilter` on the retrieved data to select sequences of behaviors such as successful and unsuccessful contacts, scanning, and miconfigurations. The resulting network flow records are then displayed with tools such as `rwcut`, `rwuniq`, `rwcount`, and `rwstats`. These tools summarize and display network flow records using specified bins in order for analysts to verify and group data for traffic characterization and behavioral analysis. Top-*N* and bottom-*N* statistics, time series, event sequence, and record-by-record displays are a few examples depicted in this analysis.

Hosts that match specific characteristics or behaviors during an analysis are then saved to named SiLK IPsets. IPsets are data structures that represent an arbitrary collection of individual addresses, and are commonly named using a behavior, characteristic, or some other descriptive attribute. For example, `webservers.set` could be a IPset file of the source IP addresses obtained from querying network flow data for flows where the source IP address responded to a SYN scan on its port 80. These binary data structures enable analysts to use the SiLK tool suite to describe network traffic and save, display, or query the hosts that match those descriptions with tools such as `rwset`, `rwsetbuild`, `rwsetcat`, or `rwfilter`.

## 2.2 Single-path Analysis: Analytics

The commands, parameters, and examples described in this chapter serve as the building blocks for analyses with the SiLK tool suite.

### 2.2.1 Get a List of Sensors With `rwsiteinfo`

The first step in a basic, single-path network analysis of the dataset described in Section 1.7 is to find out which sensor recorded the data to be analyzed and narrow down the time period for our analysis. Since routing is relatively static, data from a specific IP address generally enters or leaves through the same sensor. You need to identify the sensor that covers the affected network and figure out when this sensor recorded network flow data.

Use the `rwsiteinfo` command to view this information for the sensors on your network. `rwsiteinfo` prints SiLK configuration file information for a site, its sensors, and the traffic they collect. The `--fields` parameter is required and specifies what information is displayed. Run `rwsiteinfo` twice to do the following:

1. List the names and descriptions of all the sensors on the network. This helps to locate the sensor that covers the affected network.
2. For the sensor of interest, list the types of SiLK traffic that it carries, the number of data files stored in the SiLK repository for each type of traffic, and the start and end times for storing network flow data in the repository. This identifies the direction and type of network traffic that the sensor recorded and the time period when it was actively storing data.

Example 2.1 shows the two `rwsiteinfo` commands and their output. The results of these two calls to `rwsiteinfo` will be used in Section 2.2.2 to build a query with the `rwfilter` command to select the network flow records for our analysis.

### Determine Which Sensor Covers the Affected IP Addresses

To start, run the `rwsiteinfo` command to find the names and locations of the sensors in the network.

```
rwsiteinfo --fields=sensor,describe-sensor
```

The `--fields` parameters requests the following information:

- `sensor` displays the name of each sensor in plain text.
- `describe-sensor` displays the description of each sensor from the site configuration file (normally `silk.conf` in the root of the repository). A site's owner can specify information about the sensor configuration in this file. This gives you information (such as the sensor's location) that can help you to find which sensors recorded network traffic for the affected address block. (If the site's owner did not include this information in the site configuration file, nothing is displayed for this parameter.)

The output at the top of Example 2.1 lists the names and locations of the sensors. You need to find the sensor that covers the affected network. We are interested in traffic through the subnetwork `Div1Ext`. The sensor `S1` is associated with this subnetwork, which we will examine more closely.

### Find Traffic Types and Repository Storage Times

Once you have found the sensor of interest (`S1`), you can find out what kinds of traffic the sensor carries and when it wrote data to the SiLK repository.

```
rwsiteinfo --sensor=S1 --fields=type,repo-file-count,repo-start-date,repo-end-date
```

- `--sensor` specifies which sensor to examine. In this example, it is the name of the sensor identified via the first `rwsiteinfo` command (`S1`).
- `--fields` displays the following information in table format for sensor `S1`:

Example 2.1: Using `rwsiteinfo` to List Sensors, Display Traffic Types, and Show Repository Information

```
<1>$ rwsiteinfo --fields=sensor,describe-sensor
Sensor|Sensor-Description|
S0|Div0Ext|
S1|Div1Ext|
S2|Div0Int|
S3|Div1Int1|
S4|Div1Int2|
S5|Div1log1|
S6|Div1log2|
S7|Div1log3|
S8|Div1log4|
S9|Div1ops1|
S10|Div1ops2|
S11|Div1ops3|
S12|Div1svc|
S13|Div1dhq|
S14|Div1dmz|
S15|Div1mar|
S16|Div1med|
S17|Div1nusr|
S18|Div1mgt|
S19|Div1intel1|
S20|Div1intel2|
S21|Div1intel3|
<2>$ rwsiteinfo --sensor=S1 \
--fields=type,repo-file-count,repo-start-date,repo-end-date
Type|File-Count|Start-Date|End-Date|
in|441|2015/06/02T13:00:00|2015/06/18T18:00:00|
out|512|2015/06/02T13:00:00|2015/06/18T18:00:00|
inweb|328|2015/06/02T13:00:00|2015/06/18T18:00:00|
outweb|446|2015/06/02T13:00:00|2015/06/18T18:00:00|
innull|0|||
outnull|0|||
int2int|511|2015/06/02T13:00:00|2015/06/18T18:00:00|
ext2ext|204|2015/06/02T13:00:00|2015/06/18T18:00:00|
inicmp|0|||
outicmp|0|||
other|0|||
```

**type** – the types of enterprise network traffic that are associated with **S1**. This tells you the direction and origin of the network traffic it carries. It is also useful for splitting off data of interest (for instance, separating inbound Web traffic from other inbound traffic), which can speed up SiLK queries. (To learn more about the basic SiLK network types, see Sections 1.2.6 and 1.2.8.)

**repo-file-count** – the number of files that **S1** stored in the SiLK repository for each type of network traffic. Each file represents one hour of recorded data.

**repo-start-date** – the time and date of the oldest file that **S1** stored in the SiLK repository.

**repo-end-date** – the time and date of the most recent file that **S1** stored in the SiLK repository.

The output at the bottom of Example 2.1 lists the different types of network traffic carried by **S1**. The bulk of this traffic was recorded from 2015/06/02T13:00:00 through 2015/06/18T18:00:00. In the next step of our analysis, we will therefore retrieve network flow records from **S1** within this time period.

The output from Example 2.1 can also tell us whether **S1** recorded enough data to support a meaningful network analysis. The repository contains 441 files of inbound traffic from the ISP to the network (**in**), representing 441 hours of recorded inbound traffic to the IP addresses covered by **S1**. Similarly, the repository contains 512 hours of outbound traffic from these IP addresses to the ISP (**out**), 328 hours of inbound Web traffic (**inweb**), and 446 hours of outbound web traffic (**outweb**). This is sufficient for our analysis.

### Other Useful `rwsiteinfo` Options

Keep the following in mind when using this command:

- You must always specify parameters with `rwsiteinfo`; there is no default output.
- Enter `rwsiteinfo --fields` options in the order that you would like them to be displayed. For instance, to view the sensor description before the sensor name, specify `--fields=describe-sensor,sensor`.
- To find the classes and types supported by an installation, run `rwsiteinfo --fields=class,type,mark-defaults`. This produces three columns labeled **Class**, **Type**, and **Defaults**. The **Defaults** column shows plus signs (+) for all the types in the default class and asterisks (\*) for the default types in each class.
- The `rwsiteinfo` command supports optional parameters to control the formatting of its output (disable column spaces, change separation character, disable column headers, change field separators). It can also limit output to specific network types of interest. For these and other commonly-used parameters, see Appendix C.2. For a full list of `rwsiteinfo` options, type `rwsiteinfo --help`.

### 2.2.2 Choose Flow Records With `rwfilter`

A key step in performing a network analysis is to find and retrieve network flow records associated with the event from the SiLK repository. Use the `rwfilter` command to pull network flow records that were recorded by the sensor of interest (S1) during the time period of interest. These records will be used in subsequent steps in our analysis.

During this step in the analysis, `rwfilter` will be used to save network flow records of interest to a file. Later, we'll use `rwfilter` in conjunction with other SiLK commands to partition and explore this data.

#### About the `rwfilter` command

`rwfilter` is the most commonly used SiLK command and serves as the cornerstone for building a network analysis. It selects records from the SiLK repository, then directs the output to either files or other SiLK commands. Alternatively, `rwfilter` can select records from a pipe or file in a working directory (for instance, the output of a prior `rwfilter` command). It can optionally compute basic statistics about the flow records it reads from the repository or a file. `rwfilter` can be used on its own or in conjunction with other SiLK analysis tools, including additional invocations of `rwfilter`.

The following is a high-level view of the `rwfilter` command and its options:

```
rwfilter {selection | input} partition output
```

Specify input to `rwfilter` by using either selection or input parameters.

- *Selection* parameters read (or pull) network flow records of interest that were recorded by sensors and stored in the SiLK flow repository. They specify the attributes of records to be read from the repository, such as the sensor that recorded the data, the type of network data, the start and end dates for retrieving data, and the location of the repository.
- *Input* parameters read network flow records from pipes and/or named files in working directories containing records previously extracted from the repository or created by other means. They can be filenames (e.g., `infile.rw`) or pipe names (e.g., `stdin` or `/tmp/my.fifo`) to specify locations from which to read records. As many names as desired may be given, with both files and pipes used in the same command.

In this step of our network analysis, we will use `rwfilter`'s selection parameters to retrieve records from the SiLK repository. In future steps, we will use `rwfilter`'s input parameters to read flow records from a file or pipe.

*Partitioning* parameters create the “filter” part of `rwfilter`. These parameters specify which records pass the filter and which fail. This enables you to find and isolate network flow records that match the partitioning criteria you specify. `rwfilter` offers a variety of filtering parameters for specifying the criteria for pass/fail filtering, including time period, value ranges for packets and bytes, IP address, protocol, source and destination ports, and more.

**Hint:** An analysis will involve at least one call to `rwfilter` unless you are looking at records saved from a previous analysis. Each `rwfilter` call must include at least one partitioning parameter unless `--all-destination` is specified as an output parameter. Note that the partitioning parameter does not have to filter anything; it just needs to be present. The partitioning parameter `--protocol=0` is often used in this situation since it will not filter any records.

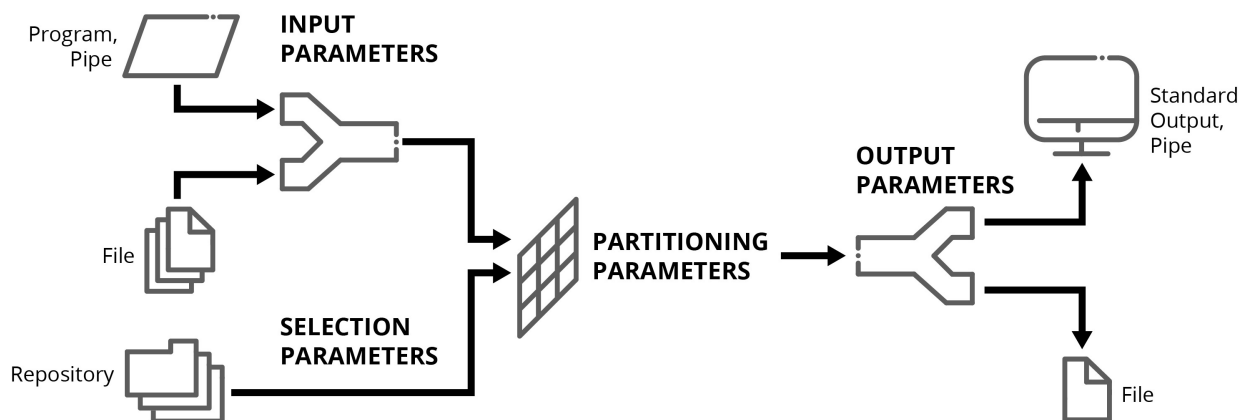
In this step of our network analysis, we will specify just one partitioning parameter: the IP address that is associated with the event. In future steps, we will specify additional partitioning parameters to identify records of interest and isolate them for further exploration with other SiLK commands.

*Output* parameters specify which group of records is returned from the call to `rwfilter`: those that “pass” the filter, those that “fail” the filter, both, or neither. These records can be written to pipes and/or named files in a working directory via the output parameters. (This also applies to statistics computed with `rwfilter`.) Each call to `rwfilter` must have at least one output parameter.

In this step of our network analysis, we will use output parameters to specify the name of the file where records are stored. In future steps, we will use pipes to direct `rwfilter` output to other SiLK commands for further investigation and processing.

Figure 2.2 shows how the `rwfilter` parameters interact.

Figure 2.2: `rwfilter` Parameter Relationships



### Retrieving network flow records and saving them to a file

Our sample single-path analysis pulls network flow records from the SiLK repository with `rwfilter`, saves them to a binary file, then examines the data in the file. This is often much faster and more efficient than pulling fresh data from the repository at every step in the analysis. `rwfilter` queries into large repositories can take a long time to run—especially if you are investigating activity over an extended period of time. To look at another group of records from the repository (for instance, from a different sensor or time period), simply run `rwfilter` again to retrieve the desired records and create additional files for analysis.

Use the `rwfilter` command as follows to pull network flow records associated with the sensor, time period, and IP address of interest to our analysis:

```
rwfilter --start=2015/06/17T14 --end=2015/06/17T14 --sensor=S1 --type=all
--any-address=192.168.70.10 --pass=flows.rw
```

- `--start` and `--end` specify the time period for retrieving records from the SiLK repository (which was found in Section 2.2.1). Times can be expressed in the form `YYYY/MM/DDThh:mm:ss.mmm`. This example uses an abbreviated time form. Since the `--start` and `--end` parameters only specify time down to the hour (e.g., `2015/06/17T14`), `rwfilter` retrieves an entire hour's worth of network flow data. One hour of traffic is pulled from the repository, starting at UTC 14:00:00.000 and ending at UTC 14:59:59.999 on June 17, 2015

**Hint:** The full parameter names are `--start-date` and `--end-date`. SiLK will recognize a parameter as long as you specify enough of its name to uniquely identify it.

- `--sensor` specifies which sensor's records to retrieve (sensor S1, which was also identified in Section 2.2.1).
- `--type` specifies the types of SiLK network traffic to retrieve. We will pull records for `all` network traffic.
- `--any-address` sets up a simple pass/fail filter for partitioning the selected network flow records. We are interested in traffic associated with the IP address 192.168.70.10. Records that match the specified IP address pass the filter; records that do not, fail it.
- `--pass` specifies the destination of the selected records that pass the filter. In this case, they are stored in the local disk file `flows.rw`.

**Hint:** Be aware that saving `rwfilter` output to a network disk file can slow down this command considerably. The speed at which records are written to the file is limited by the speed of the network. Saving to a local file is faster.

The resulting binary file, `flows.rw`, contains network flow records from the time period, sensor, traffic types, and IP address of interest. In other words, the records in this file are a snapshot of the event that we will be investigating over the course of our network analysis.

#### Example 2.2: Using `rwfilter` to Retrieve Network Flow Records From The SiLK Repository

```
<1>$ rwfilter --start=2015/06/17T14 --end=2015/06/17T14 \
--sensor=S1 --type=all --any-address=192.168.70.10 \
--pass=flows.rw
<2>$ ls -l flows.rw
-rw-r--r--. 1 analyst analyst 365935 Nov  3 14:48 flows.rw
```

### Other Useful `rwfilter` Options

Keep the following in mind when using `rwfilter`:

- Some selection parameters (such as `--sensor` and `--type`), can be used as partitioning parameters when `rwfilter` is pulling network flow records from a file or pipe. See Table C.2 for a complete list of parameters that can perform double duty for selecting and partitioning records.
- When specifying selection parameters, experienced analysts include a `--start-date` to avoid having `rwfilter` implicitly pull all records from the current day, potentially leading to inconsistent results.

- `rwfilter` partitioning parameters give analysts great flexibility in describing which flow records pass or fail the filter. Figuring out how to partition data to filter out unwanted records can be the most difficult part of using this command. Many commonly-used parameters are listed in Appendix C.3; see `rwfilter --help` for a full listing.
- Narrowing the selection of files from the repository always improves the performance of a query. On the other hand, increasing the specificity of partitioning options could improve or diminish performance. Increasing the number of partitioning parameters means more processing must be performed on each flow record. Most partitioning options involve minimal processing, but some involve considerable processing.

Generally, processing partitioning options is much less of a concern than the number of output operations, especially disk operations, and most especially network disk operations. Choosing to output both the “pass” and “fail” sets of records will involve more output operations than choosing only one set.

- The parameter `--print-filenames` lists, on the standard error file, the name of each file as `rwfilter` opens it for reading. This provides assurance that the expected files were read and indicates the command’s progress. (This is especially useful when many files are used as data sources and the command will take a long time to complete.)
- `rwfilter` can take multiple files and pipes as input. If the number of files exceeds what is convenient to put in the command line, use the `--xargs` parameter. It specifies the name of a file containing filenames from which to read flow records. This parameter also is used when another UNIX process is generating the list of input files, as in

```
find . -name '*.rw' | rwfilter --xargs=stdin...
```

### 2.2.3 View Flow Records With `rwcut`

Translating network flow records from binary format into human-readable text is a helpful part of a network analysis. Use the `rwcut` command to translate the binary network flow records selected via the `rwfilter` command as tables of ASCII text.

SiLK uses binary data to speed up queries, file manipulation, and other operations. However, these data cannot be read using any of the standard text-processing UNIX tools. `rwcut` reads SiLK flow records and translates this binary data into pipe-delimited (1) text output. You can then view the data directly in a terminal window or read it into a text-processing, graphing, or mathematical analysis tool.

**Hint:** Keep data in binary format (i.e., `*.rw` files) for as long as possible while performing an analysis. Binary SiLK network flow records are more compact and offer faster performance than the ASCII representation of these records. Use `rwcut` to inspect records or export data to other tools for further analysis.

`rwcut` can be invoked in two ways: by reading a file or by connecting it with another SiLK tool (often `rwfilter` or `rwsort`) via a pipe. When reading a file, specify the file name in the command line. The `--fields` parameter selects, reorders, and formats SiLK data fields as text and separates them in different ways.



## Displaying Flow Records

As part of the network analysis, we will use `rwcut` to take a closer look at a set of flow records from the file `flows.rw`.

```
rwcut --fields=sip,dip,sport,dport,protocol,stime --num-recs=10 flows.rw
```

- The `--fields` parameter specifies which fields in a SiLK record are shown. Field names are case-insensitive. This example displays the following fields:

`sip` – source IP address for the flow

`dip` – destination IP address for the flow

`sport` – source port for the flow

`dport` – destination port for the flow

`protocol` – transport-layer protocol for the flow

`stime` – start time of the flow, formatted as `YYYY/MM/DDThh:mm:ss.mmm`

- The `--num-recs` parameter determines how many records `rwcut` displays. In this example, up to ten records are shown (regardless of how many records are actually in the file). If the file contains no records, `rwcut` only displays the column heading for each field.
- `flows.rw` is the name of the file containing SiLK network flow records.

Example 2.3: `rwcut` for Displaying the Contents of Ten Flow Records

```
<1>$ rwcut --fields=sip,dip,sport,dport,protocol,stime \
--num-recs=10 --ipv6-policy=ignore flows.rw
  sIP|          dIP|sPort|dPort|pro|          sTime|
10.0.40.83| 192.168.70.10|53981| 8082| 6|2015/06/17T14:00:02.631|
10.0.40.20| 192.168.70.10| 53|58887| 17|2015/06/17T14:00:04.619|
10.0.40.20| 192.168.70.10| 53|55004| 17|2015/06/17T14:00:04.621|
10.0.40.83| 192.168.70.10|53982| 8082| 6|2015/06/17T14:00:12.673|
10.0.40.20| 192.168.70.10| 53|64408| 17|2015/06/17T14:00:14.685|
10.0.40.20| 192.168.70.10| 53|57734| 17|2015/06/17T14:00:14.689|
10.0.40.83| 192.168.70.10|53983| 8082| 6|2015/06/17T14:00:22.709|
10.0.40.20| 192.168.70.10| 53|63770| 17|2015/06/17T14:00:24.753|
10.0.40.20| 192.168.70.10| 53|53374| 17|2015/06/17T14:00:24.755|
10.0.40.83| 192.168.70.10|53984| 8082| 6|2015/06/17T14:00:32.741|
```

The six fields specified with the `rwcut` command are displayed in the order in which they are listed. Each field is in a separate column with its own header. The source IP addresses (`sip`) for each record vary; two addresses are shown in this example. The destination IP address (`dip`) is the same for all of these records since we only pulled records associated with that IP address.

**Hint:** Your output may contain additional spaces in the IP address field. The environment variable `SILK_IPV6_POLICY=ignore` ignores any flow record marked as IPv6, regardless of the IP addresses it contains. Only records marked as IPv4 will be printed. Setting this environment variable has the same effect as invoking `rwcut` with the `--ipv6-policy=ignore` parameter.

## Other Useful `rwcut` Options

Keep the following in mind while using the `rwcut` command:

- The `--fields` parameter selects which network flow record fields appear in `rwcut` output. Each field is associated with a number as well as a name. Table 1.1 lists the field numbers and their corresponding field names. Numbers can be specified individually or as ranges. Field names and numbers can be combined and can be listed in arbitrary order. For instance, `--fields=1-4,9,protocol` produces the same output as `--fields=sip,dip,sport,dport,stime,protocol`
- The `--fields` parameter also specifies the order in which fields are shown in output. Fields can be displayed in any order. Example 2.4 displays the output fields in this order: source IP address, source port, start time, destination IP address.

Example 2.4: `rwcut --fields` to Rearrange Output

```
<1>$ rfilter flows.rw --protocol=6 --max-pass-records=4 \
--pass=stdout | rwcut --fields=1,3,sTime,2
      sIP|sPort|          sTime|          dIP|
10.0.40.83|53981|2015/06/17T14:00:02.631| 192.168.70.10|
10.0.40.83|53982|2015/06/17T14:00:12.673| 192.168.70.10|
10.0.40.83|53983|2015/06/17T14:00:22.709| 192.168.70.10|
10.0.40.83|53984|2015/06/17T14:00:32.741| 192.168.70.10|
```

- If `--fields` is not specified, `rwcut` prints the source and destination IP address, source and destination port, protocol, packet count, byte count, TCP flags, start time, duration, end time, and the sensor name.
- The `--delimited=C` parameter changes the separator from a pipe (`|`) to any other single character, where `C` is the delimiting character. It also removes spacing between fields. This is particularly useful with `--delimited=','` which produces comma-separated-value (CSV) output for easy import into spreadsheet programs and other tools that accept CSV files. `--delimited` is the equivalent of specifying `--no-columns --no-final-delimiter --column-sep=C`.
- When output is sent to a terminal, `rwcut` (and other text-outputting tools) automatically invoke the command listed in the user's `PAGER` environment variable to paginate the output. The command given in the `SILK_PAGER` environment variable will override the user's `PAGER` environment. If `SILK_PAGER` contains the empty string, no paging will be performed. The paging program can be specified for an individual command invocation by using its `--pager` parameter.
- For a list of the most commonly-used `rwcut` parameters and options, see Appendix C.6. For a complete list of all `rwcut` parameters, enter `rwcut --help`.

### 2.2.4 Viewing File Information with `rwfileinfo`

Analyses using the SiLK tool suite can become quite complex, with several intermediate files created while isolating the behavior of interest. The `rwfileinfo` displays a variety of characteristics for each file format produced by the SiLK tool suite, which helps you to manage these files. Use this command to find out

more information about the file `flows.rw`, which contains the SiLK records associated with the IP address of interest.

For most analysts, the three most important file characteristics are the number of records in the file, the size of the file, and the SiLK commands that produced the file. Enter the following `rwfileinfo` command to view this information for `flows.rw`:

```
rwfileinfo --fields=count-records,file-size,command-lines flows.rw
```

- `--fields` specifies which SiLK file characteristics are displayed.
  - `count-records` – the total number of network flow records in a flow record file.
  - `file-size` – the size of the file in bytes.
  - `command-lines` – the commands used to generate the file. This can be very helpful when performing an analysis that involves many steps and repeated applications of commands such as `rwfilter`.
- `flows.rw` is the name of the file containing SiLK network flow records.

Output is shown in Example 2.5. `flows.rw` contains 21,864 network flow records; its size is 365,935 bytes. This gives us an idea of how much network traffic is stored there. The SiLK command that generated the file is the `rwfilter` command described in Section 2.2.2.

Example 2.5: `rwfileinfo` Displays Flow Record File Characteristics

---

```
<1>$ rwfileinfo --fields=count-records,file-size,command-lines \
  flows.rw
flows.rw:
count-records      21864
file-size          365935
command-lines
                  1  rwfilter --start=2015/06/17T14 --end=2015/06/17T14 \
--sensor=S1 --type=all --any-address=192.168.70.10 --pass=flows.rw
```

---

### Other Useful `rwfileinfo` Options

Keep the following in mind while using the `rwfileinfo` command:

- While `rwfileinfo` is generally associated with flow record files, it can also show information on sets, bags, and prefix maps (or pmaps). For more information, see Section 2.2.8, Section 4.2.4, and Section 6.2.7, respectively.
- Be sure to use the `--fields` parameter to choose which network flow record fields are displayed. If no fields are specified, `rwfileinfo` defaults to displaying a dozen fields—many of which are of no use to analysts.

- For flow record files, the record count is the number of flow records in the file. For files with variable-length records (indicated by a `record-length` of one) the field does *not* reflect the number of records; instead it is the uncompressed size (in bytes) of the data section. Notably, `count-records` does not reflect the number of addresses in an IPset file.
- Appendix C.22 lists the most commonly-used options for the `rwfileinfo` command. For a complete list of all parameters, enter `rwfileinfo --help`.

## 2.2.5 Profile Flows With `rwuniq` and `rwstats`

The next step in our network analysis is to investigate network flows to and from the IP address of interest. We will determine the most common protocols associated with these flows and find flows with low, medium, and high byte counts.

Two SiLK commands can perform these tasks.

- `rwuniq` is a general-purpose counting tool. It reads binary SiLK flow records from a file (or standard input) and counts the records, bytes, and packets for any combination of fields. `rwuniq` also groups (or *bins*) the records by a time interval specified by the analyst. In our example, this command is used to identify the hour-long time bins containing flows with low, medium, and high byte counts.
- `rwstats` provides a collection of statistical summary and counting facilities that organize and rank traffic according to various attributes. It reads binary SiLK flow records from a file (or standard input) and groups them according to a key composed of user-specified flow attributes, such as bytes, records, and packets. It then bins the records by a user-specified time interval, sums up the values of the key attributes, and sorts the bins. `rwstats` can compute statistics for each SiLK data type or for the  $n$  highest or  $n$  lowest bins. It also sums up the attribute values across all of the records it counts and displays the count for each bin as a percentage of the total. In our example, the `rwstats` command is used to identify the most commonly-used protocols associated with traffic to and from the IP address of interest.

`rwuniq` and `rwstats` overlap in their functions. Both assign flows to time bins whose length is set by the analyst. The bin size represents the length of time in seconds during which each group of records was collected, not the number of records in each bin. For each value of a key (specified by the `--fields` parameter), a bin contains counts of flows, packets, or bytes, or some other measure (specified with the `--values` parameter). `rwuniq` displays one row of output for every bin that falls within a threshold specified by the analyst. `rwstats` displays one row of output for each bin in the top  $N$  or bottom  $N$  of the total count, and computes the percentages of each data types. For a more detailed discussion of when to use each command, see [Comparing `rwstats` to `rwuniq`](#) (later in this section).

### Finding Low, Medium, and High-Byte Flows with `rwuniq`

First, use the `rwuniq` command to profile flows by byte count. It can find out how many network flow records within an hour-long period have a low byte count (between zero and 300 bytes), a medium byte count (between 300 and 100,000 bytes), or a high byte count (more than 100,000 bytes). This gives you an estimate of the volume of network activity associated with the IP address of interest.

To perform this analysis, use the `rwuniq` command in conjunction with the `rwfilter` command.

1. Run `rwfilter` on the `flows.rw` file. This file contains all traffic to and from the IP address of interest during the time period of interest; it was extracted in Section 2.2.2. Running `rwfilter` on it a second time pulls all of the records in the file with the specified byte ranges.
2. Use the Unix pipe (`|`) command to direct the resulting output to the `rwuniq` command. This command counts the number of records with each range of bytes and directs the output to a file.

```
<1> $ rwfilter flows.rw --bytes=0-300 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type --values=records --sort-output \
> low-byte.txt

<2> $ rwfilter flows.rw --bytes=300-100000 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type --values=records --sort-output \
> medium-byte.txt

<3> $ rwfilter flows.rw --bytes=100000- --pass=stdout --values=records \
| rwuniq --bin-time=3600 --fields=stime,type --sort-output \
> high-byte.txt
```

Parameters for the `rwfilter` command include the following:

- `flows.rw` contains the network flow records of interest.
- `--bytes` specifies the range of byte counts for selecting records. Ranges are specified using a dash (e.g., `0-300` selects all flows with byte counts between zero and 300). To specify an open-ended range, do not include an upper bound on the range (e.g., `100000-` selects all flows with byte counts greater than 100000).
- `--pass=stdout` sends all records that pass the filter to standard output.

Parameters for the `rwuniq` command include the following:

- `flows.rw` is the name of the file containing SiLK network flow records.
- `--bin-time=3600` defines a time bin that is one hour (3600 seconds) long.
- `--fields=stime,type` specifies the fields to use as keys for counting network flows. This parameter is required. We are looking at the values for `stime` (start time for the flow) and `type` (network flow type).
- `--values=records` counts the number of records that passed the `rwfilter` command.
- `--sort-output` sorts the output of the `rwuniq` command in numerical order according to the value (or values) of the key specified via the `--fields` parameter.
- The shell command `>` directs the output of `rwuniq` into the file `low-byte.txt`.

**Hint:** We could have saved the `rwfilter` output to a file and run `rwuniq` on that file instead of using the UNIX pipe (`|`) command to send the output directly to the `rwuniq` command. However, one problem with generating such temporary files is that they slow down the analysis. The `rwfilter` command would have written all the data to disk, and then the subsequent `rwuniq` command would have read the data back from disk. Using UNIX pipes to pass records from one process to another skips the time-consuming steps of writing and reading data, speeding this up considerably. The SiLK tools can operate concurrently, using memory (when possible) to pass data between them. Setting up an unnamed pipe between processes is described in Appendix B.2.

Example 2.6 shows the output from this series of SiLK commands.

Example 2.6: Characterizing flow byte counts with `rwuniq`

```
<1>$ rfilter flows.rw --bytes=0-300 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >low-byte.txt
<2>$ cat low-byte.txt
      sTime|   type|  Records|
2015/06/17T14:00:00|   in|    1449|
2015/06/17T14:00:00|   out|    1500|
2015/06/17T14:00:00| inweb|     12|
2015/06/17T14:00:00| outweb|    8339|
<3>$ rfilter flows.rw --bytes=300-100000 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >medium-byte.txt
<4>$ cat medium-byte.txt
      sTime|   type|  Records|
2015/06/17T14:00:00|   in|     66|
2015/06/17T14:00:00|   out|     96|
2015/06/17T14:00:00| inweb|    346|
2015/06/17T14:00:00| outweb|   10051|
<5>$ rfilter flows.rw --bytes=100000- --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >high-byte.txt
<6>$ cat high-byte.txt
      sTime|   type|  Records|
2015/06/17T14:00:00|   out|     3|
2015/06/17T14:00:00| outweb|     2|
```

### Other useful `rwuniq` options

- The `--value` parameter specifies which flow attributes are counted for a time bin. In addition to counting bytes, `rwuniq` can count records, packets, and source and destination IP addresses.
- Flow records need not be sorted before being passed to `rwuniq`. If the records are sorted in the same order as indicated by the `--fields` parameter to `rwuniq`, using the `--presorted-input` parameter may reduce memory requirements for `rwuniq`.
- For a list of the most commonly-used `rwuniq` parameters, see Appendix C.8. For a complete list of all `rwuniq` parameters, enter `rwuniq --help`.

### Finding the Most Commonly-Used Protocols With `rwstats`

Another way to characterize network flows is by protocol usage. By looking at the most commonly-used protocols, we can get a sense of what types of traffic the network carries.

Use the `rwstats` command to identify the 10 most common protocols associated with traffic into and out of the IP address of interest. `rwstats` groups records into time bins by field (or fields), similar to `rwuniq`. However, `rwstats` can list the top  $N$  or bottom  $N$  bins and compute summary percentages for each item.

```
rwstats --fields=protocol --count=10 flows.rw
```

- `--fields=protocol` counts the records that carry traffic with each protocol.
- `--count=10` computes statistics for the 10 bins with the most common protocols.
- `flows.rw` is the name of the file containing SiLK network flow records.

Example 2.7 shows the output from this command.

Example 2.7: Finding the top protocols with `rwstats`

```
<1>$ rwstats --fields=protocol --count=10 flows.rw
INPUT: 21864 Records for 3 Bins and 21864 Total Records
OUTPUT: Top 10 Bins by Records
prol  Records| %Records|  cumul_%|
 6|      18854| 86.233077| 86.233077|
17|      2909| 13.304976| 99.538053|
 1|       101|  0.461947|100.000000|
```

Notice that the output lists just three protocols, not ten. This is because only three protocols were used during the time period of interest. `rwstats` also computes the number of records for each protocol and summarizes the percentage of traffic for each protocol.

- Protocol 6 (Transmission Control Protocol, or TCP) makes up approximately 86% of the traffic; it is used by popular applications such as the World Wide Web, email, remote administration, and file transfer programs.
- Protocol 17 (User Datagram Protocol, or UDP) makes up approximately 13% of the traffic; it is used by the Domain Name System (DNS), the Routing Information Protocol (RIP), the Simple Network Management Protocol (SNMP), and the Dynamic Host Configuration Protocol (DHCP). Voice and video is also transmitted using UDP.
- Protocol 1 (Internet Control Message Protocol, or ICMP) makes up less than 1% of the traffic; it is used by network devices (such as routers) to transmit error messages and other information.

### Other useful `rwstats` options

- Each call to `rwstats` *must* include exactly one of the following:
  - a key containing one or more fields via the `--fields` parameter and an option to determine the number of key values to show via `--count` (shown in Example 2.7), `--percentage`, or `--threshold`
  - one of the summary parameters (`--overall-stats` or `--detail-proto-stats`)
- For a list of the most commonly-used `rwstats` parameters, see Appendix C.4. For a complete list of all `rwstats` parameters, enter `rwstats --help`.

### Comparing `rwstats` to `rwuniq`

`rwstats` in top or bottom mode and `rwuniq` have much in common, especially since SiLK version 3.0.0. An analyst can perform many tasks with either tool. Some guidelines follow for choosing the tool that best suits a task. Generally speaking, `rwstats` is the workhorse data description tool, but `rwuniq` does have some features that are absent from `rwstats`.

- Like `rwcount`, `rwstats` and `rwuniq` assign flows to bins. For each value of a key, specified by the tool with the `--fields` parameter, a bin summarizes counts of flows, packets, or bytes, or some other measure determined by the analyst with the `--values` parameter. `rwuniq` displays one row of output for every bin except those not achieving optional thresholds specified by the analyst. `rwstats` displays one row of output for each bin in the top  $N$  or bottom  $N$ , where  $N$  is determined directly by the `--count` parameter or indirectly by the `--threshold` or `--percentage` parameters.
- If `rwstats` or `rwuniq` is initiated with multiple counts in the `--values` parameter, the first count is the primary count. `rwstats` can apply a threshold only to the primary count, while `rwuniq` can apply thresholds to any or several counts.
- For display of all bins, `rwuniq` is easiest to use. However, a similar result can be obtained with `rwstats --threshold=1`. `rwstats` will run more slowly than `rwuniq` because it must sort the bins by their summary values.
- `rwstats` always sorts bins by their primary count. `rwuniq` optionally sorts bins by their key.
- `rwstats` normally displays the percentage of all input traffic accounted for in a bin, as well as the cumulative percentage for all the bins displayed so far. This output can be suppressed with `--no-percents` to be more like `rwuniq` or when the primary count is not bytes, packets, or records.
- `rwuniq` has two counts that are not available with `rwstats`: `sTime-Earliest` and `eTime-Latest`.
- Network traffic frequently can be described as exponential, either increasing or decreasing. `rwstats` is good for looking at the main part of the exponential curve, or the tail of the curve, depending on which is more interesting. `rwuniq` provides more control of multi-dimensional data slicing, since its thresholds can specify both a lower bound and an upper bound. `rwuniq` will be better at analyzing multi-modal distributions that are commonly found when the x-axis represents time.

### 2.2.6 Characterize Traffic by Time Period With `rwcount`

A typical network analysis will examine network traffic by time period to see how it varies throughout the event of interest. Unusual volumes of traffic, changes in byte and packet counts, and other deviations from normal activity can help you to figure out what is causing the event to occur.

The `rwcount` command captures network activity that occurs during the time interval (or bin) that you specify. It counts the number of records, bytes, and packets for flows occurring during a bin's assigned time period. You can then view these counts in a terminal window or graph them in a plotting package such as gnuplot, a spreadsheet package such as Microsoft Excel, or another analysis tool.

Our analysis will examine network traffic to and from the target IP address during the time period of interest. We will use `rwcount` to show this activity in ten-minute time bins.

```
rwcount --bin-size=600 flows.rw
```



- `--bin-size` specifies the time bin in seconds. In this example, the time bin is 600 seconds or ten minutes.
- `flows.rw` contains the network flow records of interest.

Example 2.8 shows the output from this command. It counts the flow volume information gleaned from the `flows.rw` file by ten-minute bins.

Example 2.8: Counting Bytes, Packets and Flows with Respect to Time

Date	Records	Bytes	Packets
2015/06/17T14:00:00	466.00	798757.00	3423.00
2015/06/17T14:10:00	394.00	104668.00	1622.00
2015/06/17T14:20:00	382.43	104159.18	1621.86
2015/06/17T14:30:00	393.57	107100.82	1670.14
2015/06/17T14:40:00	9335.01	15559931.61	191709.67
2015/06/17T14:50:00	10885.11	16541697.17	187619.55
2015/06/17T15:00:00	7.70	75830.56	897.45
2015/06/17T15:10:00	0.17	21466.66	383.33

By default, `rwcount` produces the table format shown in Example 2.8.

- The first column is the timestamp for the earliest moment in the bin.
- The net three columns show the number of flow records, bytes, and packets counted in the bin.

### Examining Bytes, Packets, and Flows

Counting by bytes, packets, and flows can reveal different traffic characteristics. As noted at the beginning of this manual, the majority of traffic crossing wide area networks has very low packet counts. However, this traffic, by virtue of being so small, does not make up a large volume of bytes crossing the enterprise network. Certain activities, such as scanning and worm propagation, are more visible when considering packets, flows, and various filtering criteria for flow records.

The traffic into and out of the IP address of interest (captured in the file `flows.rw`) jumps significantly during the ten-minute time bins `2015/06/17T14:40:00` and `2015/06/17T14:50:00`. Byte, packet, and record counts all rise during this 20-minute time period.

### Examining Traffic Over a Period of Time

`rwcount` is used frequently to provide graphs showing activity over long periods of time, giving a visual representation of shifts in network traffic. Count data can be read by most plotting (graphing) applications.

The data from Example 2.8 is plotted using Microsoft Excel in Figure 2.3. The traffic spike that we saw in the tabular data shows up clearly in the plots on the left-hand side of this figure.

For a more detailed look at network activity during this time period, we can change the `--bin-size` from 600 seconds (ten minutes) to 60 seconds (one minute).

```
rwcount --bin-size=60 flows.rw
```

Plots of this data are shown on the right-hand side of Figure 2.3. Looking at the data on a minute-by-minute basis shows the variation in data flows during this event.

**Hint:** Whether you use a larger bin size or smaller bin size depends on your data. Smaller bin sizes provide more data points to capture subtleties in traffic. If the bin size is too small, however, it becomes harder to spot trends in the data. Larger bin sizes make it easier to spot regular traffic patterns. If the bin size is too large, however, there will not be enough resolution in the data to see what is happening on your network at a given point in time.

Figure 2.3: Displaying `rwcount` Output Using 10-Minute and 1-Minute Bins



### Other Useful `rwcount` Information

Keep the following in mind when using `rwcount`.

- The default bin size is 30 seconds.
- Bin counts that have zero flows, packets, and bytes can be suppressed by the `--skip-zeroes` option to reduce the length of the listing. However, do not skip rows with zero flows if the output is being passed to a plotting program; if they are, those data points will not be plotted.

### 2.2.7 Sort Flow Records With `rwsort`

Sorting flow records can help you to organize them according to protocol, IP address, start time, and other attributes. Use the `rwsort` command to sort binary flow records according to the value of the field(s) you select.

`rwsort` is a high-speed sorting tool for SiLK flow records. It reads binary SiLK flow records from a file (or standard input) and outputs the same records in a user-specified order. The output records are also in binary (non-text) format and are not human-readable without interpretation by another SiLK tool such as `rwcut`. `rwsort` is faster than the standard UNIX `sort` command, handles flow record fields directly with understanding of the fields' types, and is capable of handling very large numbers of SiLK flow records provided sufficient memory and storage are available.

The following example sorts the network flow records in `flows.rw` by byte count, destination IP, and protocol from the highest to the lowest value in each field, then displays the first ten records.

```
rwsort flows.rw --fields=dip,protocol,bytes --reverse
| rwcut --fields=dip,protocol,bytes,stime --num-recs=10
```

- `--fields` specifies the sort order. It identifies the fields that are used as sort keys and specifies their precedence. In this example, `rwsort` first sorts the records by destination IP address (`dip`), then protocol (`protocol`), then byte count (`bytes`).
- By default, `rwsort` sorts from the lowest to the highest values of each sort key. `--reverse` sorts the records from the highest to the lowest values.
- The file `flows.rw` contains the SiLK record files to be sorted.
- The Unix pipe command (`|`) sends the output of the `rwsort` command to the `rwcut` command.
- The `rwcut` command and its parameters are described in Section 2.2.3.

Example 2.9 shows the results of this command. The records are first sorted from the highest destination IP address to the lowest. They are then sorted according to their protocols, then their sizes in bytes. This gives you an idea of the volume and types of traffic associated with the destination IPs.

### Behavioral Analysis with `rwsort`, `rwcut`, and `rwfilter`

A behavioral analysis of protocol activity relies heavily on basic `rwcut` and `rwfilter` parameters. The analysis requires the analyst to have a thorough understanding of how protocols are meant to function. Some concept of baseline activity for a protocol on the network is needed for comparison.

Example 2.9: Sorting by Destination IP Address, Protocol, and Byte Count

```

<1>$ rwsort flows.rw --fields=dip,protocol,bytes --reverse \
| rwcut --fields=dip,protocol,bytes,stime --num-recs=10
      dIP|pro|      bytes|      sTime|
216.207.68.32| 6|      960|2015/06/17T14:53:15.707|
216.207.68.32| 6|      960|2015/06/17T14:54:30.604|
216.207.68.32| 6|      120|2015/06/17T14:54:14.405|
216.207.68.32| 6|      120|2015/06/17T14:55:29.333|
209.66.102.50| 6|      960|2015/06/17T14:55:26.586|
209.66.102.50| 6|      960|2015/06/17T14:56:42.465|
209.66.102.50| 6|      120|2015/06/17T14:57:41.186|
209.66.102.50| 6|      120|2015/06/17T14:56:25.264|
208.206.41.61| 6|      960|2015/06/17T14:58:20.666|
208.206.41.61| 6|      960|2015/06/17T14:46:17.427|

```

To monitor the behavior of protocols, first take a sample of a particular protocol. Use `rwsort --fields=sTime`, and convert the results to ASCII text with `rwcut`. To produce byte and packet fields only, try `rwcut` with `--fields=bytes` and `--fields=packets`. Then, perform the UNIX commands `sort` and `uniq -c`.

Cutting in this manner (sorting by field or displaying select fields) can answer a number of questions:

1. Is there a standard bytes-per-packet ratio?
2. Do any bytes-per-packet ratios fall outside the baseline?
3. Do any sessions' byte counts, packet counts, or other fields fall outside the norm?

There are many such questions to ask, but keep the focus of exploration on the behavior being examined. Chasing down weird cases is tempting but can add little to your understanding of general network behavior.

### Other Useful `rwsort` Information

Keep the following in mind when using `rwsort`.

- Sort keys can be specified by field numbers as well as field names; see Table 1.1 for a complete list.
- Sort keys can be specified in any order. For example, `--fields=1,3` results in flow records being sorted by source IP address (1) and by source port (3) for each source IP address. `--fields=3,1` results in flow records being sorted by source port and by source IP address for each source port. (Since flow records are not always entered into the repository in the order in which they were initiated, analyses often involve sorting by start time at some point.)
- `rwsort` can also be used to sort multiple SiLK record files. If the flow records in the input files are already ordered in each file, using the `--presorted-input` parameter can improve efficiency significantly by just merging the files.
- If `rwsort` is processing large input files, disk space in the default temporary system space may be insufficient or not located on the fastest storage available. To use an alternate space, specify the `--temp-directory` parameter with an argument specifying the alternate space. This may also improve data privacy by specifying local, private storage instead of shared storage.

## 2.2.8 Use IPsets to Gather IP Addresses

Up to this point, our single-path analysis has focused on selecting, storing, and examining flow records. However, another common goal of single-path analysis is to compile lists of IP addresses that exhibit criteria of interest to the analyst. This section will continue our analysis by gathering and summarizing single-path criteria using named sets of IP addresses, or *IPsets*.

### Create IPsets With `rwset` and `rwsetbuild`

`rwset` and `rwsetbuild` are two SiLK tools for creating sets of IP addresses (IPsets). `rwset` creates sets from flow records. `rwsetbuild` creates them from lists of IP addresses in text files. Expanding on the profiling in Section 2.2.5, `rwfilter` can be used to profile network flows by bytes. When combined, `rwset` and `rwfilter` summarize the IP addresses that exhibit byte-threshold profiles to files with descriptive names.

```
rwfilter flows.rw --bytes=0-300 --pass=stdout \  
| rwset --any-file=low-byte.set
```

Parameters for the `rwfilter` command include the following:

- `flows.rw` contains the network flow records of interest.
- `--bytes=0-300` specifies the range of byte counts for selecting records (0-300 for this example).
- `--pass=stdout` sends all records that pass the filter to standard output.

Parameters for the `rwset` command include the following:

- `--any-file=low-bytes.set` specifies source and destination IP addresses from flow records with a range of 0-300 bytes to the IPset file `low-bytes.set`. Because `--any-file` was used above, the IPset file will include the IP address itself as well as any IP addresses that communicated with it.

Example 2.10 shows the output from this series of `rwfilter` and `rwset` commands.

Example 2.10: Using `rwset` to Gather IP Addresses

---

```
<1>$ rwfilter flows.rw --bytes=0-300 --pass=stdout \  
| rwset --any-file=low-byte.set  
<2>$ file low-byte.set  
low-byte.set: SiLK, IPSET v2, Little Endian, LZ0 compression
```

---

Analysis requiring defined IP addresses should use the `rwsetbuild` tool. `rwsetbuild` creates SiLK IPsets from textual input, including canonical IP addresses, CIDR notation, and IP ranges. This approach is useful for creating whitelists and blacklists of IP addresses that may reside in network flow records presently or in the future.

```
rwsetbuild --ip-ranges servers.txt servers.set
```

Parameters for the `rwsetbuild` command include the following:

- `--ip-ranges` specifies allowing the textual input file to contain IP ranges.
- `servers.txt` specifies the textual input file name.
- `servers.set` specifies the binary IPset output file name.

Example 2.11 shows the output from the `rwsetbuild` command.

Example 2.11: Using `rwsetbuild` to Gather IP Addresses

---

```
<1>$ cat servers.txt
# Text file of servers
192.168.2.1 # Single
192.168.3.0/24 # CIDR
192.168.4.1-192.168.4.128 # IP range
<2>$ rwsetbuild --ip-ranges servers.txt servers.set
<3>$ file servers.set
servers.set: SiLK, IPSET v2, Little Endian, LZ0 compression
```

---

### Other Useful `rwset` and `rwsetbuild` Options

Keep the following in mind while using the `rwset` command:

- `rwset` can assign IP addresses to IPsets by source, destination, and both source and destination simultaneously.
- `rwset` and `rwsetbuild` can read input from files on disk or standard input (`stdin`).
- `rwsetbuild` supports SiLK IP address wildcard notation (10.x.1-2.4,5). This notation is not supported when the `--ip-ranges` switch is specified.
- Appendix C.14 lists the most commonly-used options for the `rwset` command. For a complete list of all parameters, enter `rwset --help` or `rwsetbuild --help`.

### Display IP Addresses, Counts, and Network Information With `rwsetcat`

Single-path analysis often requires the IP addresses in an IPset to be counted and displayed. This gives you an opportunity to inspect the IP addresses that met specified analytic criteria, such as behavior and network topology. `rwsetcat` can display the IP addresses in an IPset, count the number of IP addresses, display information about the network, and show minimum and maximum IP addresses as well as other summary data for the IPset.

We will continue our analysis of the low byte IP addresses from Section 2.2.8 by counting, listing, and computing summary statistics for the IP addresses in the `low-bytes.set` IPset file.

To count the number of IP addresses that exhibited 0-300 byte flow records with the IP address 192.168.70.10, enter the following command:

```
rwsetcat --count-ips low-byte.set
```

Parameters for the `rwsetcat` command include the following:

- `--count-ips` specifies counting the number of IP addresses.
- `low-byte.set` specifies the binary IPset file for counting (containing IP addresses that exhibited 0-300 byte flow records with 192.168.70.10.)

Example 2.12 shows the count of IP addresses contained in `low-byte.set`.

Example 2.12: Using `rwsetcat` to Count Gathered IP Addresses

```
<1>$ rwsetcat --count-ips low-byte.set
574
```

Although general counting is helpful, an analysis commonly requires additional context regarding the networks and hosts contained in the IPset. `rwsetcat` prints analyst-specified subnet ranges and the number of hosts in each subnet.

To summarize the /24 networks contained in `low-byte.set`:

```
rwsetcat --network-structure=24 low-byte.set
```

Parameters for the `rwsetcat` command include the following:

- `--network-structure` groups IP addresses by specified structure and prints the number of hosts
- `low-byte.set` specifies the binary IPset file for counting (containing IP addresses that exhibited 0-300 byte flow records with 192.168.70.10.)

Example 2.13 shows the first four /24 networks contained in `low-byte.set` and their respective host counts.

Example 2.13: Using `rwsetcat` to Print Networks and Host Counts

```
<1>$ rwsetcat --network-structure=24 low-byte.set | head -n 4
 4.2.0.0/24| 1
 6.7.1.0/24| 1
 8.1.7.0/24| 1
10.0.20.0/24| 1
```

Complete statistical summaries are also common during an analysis and can be printed with `rwsetcat`. Our previous /24 summary only prints the specified CIDR range and would require iterative commands to determine multiple CIDR network ranges that may be contained in an IPset. Therefore, `rwsetcat` provides the `--print-statistics` switch for full statistical summaries of an IPset.

```
rwsetcat --print-statistics low-byte.set
```

Parameters for the `rwsetcat` command include the following:

- `--print-statistics` specifies printing a statistical summary of IP addresses contained in an IPset.
- `low-byte.set` specifies the binary IPset file for counting (containing IP addresses that exhibited 0-300 byte flow records with 192.168.70.10.)

Example 2.14 shows the statistical summary of IP addresses in the `low-byte.set` file.

Example 2.14: Using `rwsetcat` to Print IP Address Statistical Summaries

```
<1>$ rwsetcat --print-statistics low-byte.set
Network Summary
  minimumIP =          4.2.0.58
  maximumIP =       216.207.68.32
      574 hosts (/32s),    0.000013% of 2^32
      87 occupied /8s,    33.984375% of 2^8
      381 occupied /16s,   0.581360% of 2^16
      521 occupied /24s,   0.003105% of 2^24
      551 occupied /27s,   0.000411% of 2^27
```

### Other Useful `rwsetcat` Options

Keep the following in mind while using the `rwsetcat` command:

- `rwsetcat` can print CIDR blocks without specifying a desired network mask. `rwsetcat` will group sequential IPs into the largest possible CIDR block and prints individual IP addresses. This switch cannot be combined with the `--network-structure` switch.
- The `--network-structure` switch supports multiple CIDR masks for a single command execution.
- Appendix C.15 lists the most commonly-used options for the `rwsetcat` command. For a complete list of all parameters, enter `rwsetcat --help`.

## 2.2.9 Resolve IP Addresses to Domain Names With `rwresolve`

Use the `rwresolve` command to find the host names associated with the IP addresses of interest to our network analysis. This command performs a reverse Domain Name Service (DNS) lookup on a list of IP addresses to retrieve their host names. If the lookup is successful, it prints the name of the host; if not, it prints the IP address. If an IP address resolves to multiple host names, it print the first one found. The result is a human-readable list of host names that is useful for further investigation and analysis.

`rwresolve` takes delimited text as input, not binary flow records. It is designed for use with the `rwcut` command, although it can be used with any SiLK tool that produces delimited text.

**Hint:** Since performing reverse DNS lookups is a time-consuming process, we strongly recommend that you use `rwresolve` only on small datasets.

```
rwcut --fields=1,1 flows.rw | rwresolve --ip-field=2
```



This command first uses the `rwcut` command to generate a list of IP addresses (`--fields=1,1`). It redirects the resulting output to the `rwresolve` command, which looks up the host names associated with the IP addresses.



## Chapter 3

# Case Studies: Basic Single-path Analysis

The previous chapter introduced the process of single-path analysis and covered some of the commands that are used in such analyses. This chapter walks through several detailed cases that serve as examples of single-path analyses.

Upon completion of this chapter you will be able to

- describe a sequence of steps that analysts may use in approaching a task
- apply those steps to several tasks relevant to network traffic
- use SiLK tools to automate the analysis

The case studies in this chapter use the FCCX dataset described in Section [1.7](#).

### 3.1 Profile Traffic Around an Event

One view of a network security event is that some specific activity occurs on a particular host at an identified time. In terms of the analysis in this handbook, a host is indicated by its IP address, and time is, at first consideration, associated with a given hour. With this as a starting point, the analyst needs to develop a high-level assessment of possible changes in behavior which then provides a guide to more detailed follow-on assessments. The end goal is to answer some basic questions about the event:

- Did the event impact network performance or services?
- Was the impact sufficient to warrant dedicating resources to respond?
- Was the event malicious?
- Did it demonstrate weaknesses that could enable malicious activity?
- Which entities were involved (both internal and external)?

Frequently, trying to answer such questions in detail involves too much effort. The alternative is to proceed in a staged manner, and at each stage determine if analysis should proceed further. This section describes an initial high-level analysis that can be done rapidly. It helps you to determine whether there could have been some impact from the event, along with a rough feel as to the magnitude of that impact.

Working from the target IP address and the time frame as a start, there are several possible approaches to gaining a high level indication of impact from the event:

- **traffic**—look at traffic on the targeted network and search for shifts in the size and frequency of contacts involving the target, measuring before and after the event
- **response time**—look at the overall response time for service requests (the average interval between request and response) into the network, then determine if it has increased during and following the event
- **contact rate**—look at the relative rate of contact with services (indicated by port and protocol) on the targeted network, searching for shifts in the contact rate and the size of traffic on those services
- **hosts**—look at the set of hosts in contact with the target, and determine if it has shifted unusually during and after the event.

This section will explore the first of these alternatives: looking at traffic shifts. (The other alternatives may be useful to analysts, but are not covered here.)

### 3.1.1 Examining Shifts in Traffic

We can apply the Formulate-Model-Analyze steps in the SiLK workflow to perform a single-path analysis that looks at shifts in network traffic around the event. First, we will filter for network flow records associated with the targeted host around the time of the event (the Formulate step from Section 1.5). This involves pulling records from the appropriate parts of the repository and isolating those that involve the targeted host. This set of records can then be divided into bins according to volumes, and then into counts for each bin before, during, and after the event (the Model step). Finally, the counts can then be interpreted to assess the potential impact of the event (the Analyze step).

#### Filter Traffic Around the Event

The filter portion of the analysis is structured as a query using the `rwfilter` tool. The appropriate parts of the repository are determined by date-hour (using the parameters `--start` and `--end`) and type (using the `--type` parameter). The association with the target host is indicated by the host's IP address (using the `--any-address` parameter). The filtered records are then stored in a file (using the `--pass` parameter) for later parts of the analysis.

#### Summarize Records

Once the records are pulled via the query, they are summarized by using `rwuniq`. The goal is to filter out the appropriate group of flows to summarize. We will create volume-based groups at low, medium, and high values for both byte volume and flow duration.

For each group, the analysis uses another call to `rwfilter` to pull records from the file generated previously. It extracts those with the volume measure for each group (using either `--bytes` or `--duration`). The output goes to `rwuniq` to count the records. Separate counts are generated for each hour and type of flow record (using `--fields=stime,type` and `--bin-time=3600`). The number of records in each group are counted (using `--values=records`). This provides a high-level view of the variation in activity from an hour before the event to an hour after it.

### 3.1.2 How to Profile Traffic

The resulting set of commands for this analysis are shown in Example 3.1. Command 1 is the initial query. It filters records from the repository that are associated with a specific IP address and saves them to the file `traffic.rw` in the local directory. Commands 2 through 7 are the processing steps to summarize each group of records. They produce text files with hourly counts for each type in each group.

Example 3.1: Using `rwfilter` and `rwuniq` to Profile Traffic Around an Event

---

```
<1>$ rwfilter --start=2015/06/17T13 --end=2015/06/17T15 \
  --sensor=S1 --type=in,inweb,out,outweb \
  --any-address=192.168.70.10 --pass=traffic.rw
<2>$ rwfilter traffic.rw --bytes=0-300 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >low-byte.txt
<3>$ rwfilter traffic.rw --bytes=301-100000 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >med-byte.txt
<4>$ rwfilter traffic.rw --bytes=100001- --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >high-byte.txt
<5>$ rwfilter traffic.rw --duration=0-60 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >short-duration.txt
<6>$ rwfilter traffic.rw --duration=61-120 --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >med-duration.txt
<7>$ rwfilter traffic.rw --duration=121- --pass=stdout \
| rwuniq --bin-time=3600 --fields=stime,type \
  --values=records --sort-output >long-duration.txt
```

---

The results of these commands are collated in Example 3.2. The counts show that there was a marked increase in low-to-medium byte and short-to-medium duration web traffic during and after the event. There was no corresponding increase in high byte or long duration traffic. Based on this, the analyst may start to focus to look for what common factors exist in the increased traffic. The goal is to build awareness of the impact of the event in a way that helps responders to deal with that impact.

Example 3.2: Collated Profile of Traffic Around an Event

stime	type	sbyte	mbyte	hbyte	sdur	mdur	ldur
2015/06/17T13:00:00	in	720	160		880		
2015/06/17T13:00:00	inweb	5	352		357		
2015/06/17T13:00:00	out	764	192	9	960	5	
2015/06/17T13:00:00	outweb	1	400		401		
2015/06/17T14:00:00	in	1449	66		1515		
2015/06/17T14:00:00	inweb	12	346		358		
2015/06/17T14:00:00	out	1500	96	3	1595	1	1
2015/06/17T14:00:00	outweb	8339	10051	2	18382	9	1
2015/06/17T15:00:00	in	2528	550		3077	1	
2015/06/17T15:00:00	inweb	14	345		359		
2015/06/17T15:00:00	out	2520	558	5	3076	3	3
2015/06/17T15:00:00	outweb	10309	11072	7	21366	12	9

## 3.2 Generate Top $N$ Lists

Filtering flow records by time, sensor, type, and volume characteristics often produces groups that contain flow records of interest. However, these groups also contain extraneous flows that produce noise, which makes it more difficult to spot patterns in the data.

One strategy for removing these extraneous flows is to identify the largest sub-groups, validate each sub-group, and either set it aside or include it in the collection of flows of interest. The sub-groups are identified by a combination of flow characteristics, such as the IP address of the source, the TCP flags present in the flow, or the network service involved. The contribution of each sub-group is measured by the total bytes or packets per sub-group, the number of records per sub-group, the number of distinct values present for some field in the flow record, or another summary statistic.

The overall process of pulling a collection of flow records and then removing flows not related to the analysis is sometimes referred to as *top-down analysis*. There is also a *bottom-up analysis* that involves starting with a minimal set of records that are of interest, then basing further queries to isolate more records of interest based on the field values in the minimal set.

### 3.2.1 Using `rwstats` to Create Top $N$ Lists

To identify the identity and relative size of the sub-groups, use the `rwstats` command. It explicitly includes parameters to limit output to the largest contributors and describes the contribution of those categories to the overall flow collection<sup>5</sup>. Without `rwstats`, an analyst could load flow records into a spreadsheet, then generate a pivot table to identify the most common characteristics. `rwstats` is much faster and easier to use than a spreadsheet. It deals with high numbers of flow records very efficiently in terms of storage and memory usage.

<sup>5</sup>`rwstats` has further functions that facilitate describing collections of flows statistically, which are described in Appendix C.4

## Removing Unwanted Flows

In generating top- $N$  lists, the Formulate stage involves eliminating flow records for network activity that is not of interest. This activity, sometimes referred to as network chaff, may include connections with not enough data exchanged to be significant, those involving services that are not important for an analysis, or, in general, anything that could obfuscate the results by contaminating the flow records retrieved for the event.

In Example 3.3, command 1 looks at the first few flows in the `traffic.rw` file generated in Example 3.1. The sequence of flows shown are DNS queries. There is not enough information at the flow level to indicate whether they are relevant to the event that occurred.

Command 2 in the example uses a new call to `rwfilter` to exclude the unneeded flows, both saving a copy to a new flow file and sending it to `rwcut` for further examination. The examination shows that the host at 192.168.70.10 is doing a lot of communication on TCP port 8082, associated with a file management utility known as Utilistore, and the larger-volume flows appear to be associated with the host at 10.0.40.83.

Command 3 queries the flow repository to look for flows showing communication on this port with at least 150 average bytes per packet across the full data set. The results of this query are then passed to `rwuniq` to profile all of the locations to which data has been sent. The results of this profile show three hosts receiving this traffic, including both 192.168.70.10 and another host at 192.168.200.10.

Command 4 uses a further call to `rwfilter` to pull all the traffic associated with the newly located addresses 10.0.40.83 and 192.168.70.10. The addresses appear twice in the `rwfilter` call in order to specify that both the source and destination are constrained to these addresses. The results are then stored in the file `traffic2.rw` to be examined further.

Note that the addresses are specified with `--scidr` (Source CIDR block) and `--dcidr` (Destination CIDR block) instead of `--saddress` (source address), `--daddress` (destination address), or `--any-address` (both source and destination addresses). The `--scidr` and `--dcidr` parameters accept comma-separated lists of addresses in CIDR notation, whereas the other parameters accept only a single address. The `/32` CIDR notation specifies a single address and thus permits us to use a list of addresses.

## Summarizing Destination Port Usage By Records and Bytes

After we query and filter the flow records to isolate those of interest, we can calculate values to clarify our understanding of these data (the Model step). We could use either `rwuniq` or `rwstats` to understand the contributors to these data, which fed into the filtering process. `rwstats` allows for more explicit limits on the number of bins that are displayed. In contrast, `rwuniq` shows all of the bins for the input dataset. `rwstats` also shows the percentage contribution to the overall input of each bin and cumulatively across bins. These limits are often expressed as a count of bins, but they can also be expressed in terms of percentage contribution or a threshold on the count. The percentages are calculated based only on the first value specified for the bin. This allows `rwstats` to be used flexibly to profile the contributors to the data.

In Example 3.3, command 5 uses `rwstats` to profile the records in `traffic2.rw`, looking at the destination port utilization in these data by the flow count (as a rough measure for how often communication takes place), with bytes also calculated as a supplementary value.

In the results, the largest contributor accounts for very close to half of the data. This is not surprising since the analyst used this port to identify these hosts as being of interest during the filtering process. Three of the other ports shown are ephemeral ports (officially, ports numbering 49,152 or more, although some Linux versions use 32,768 to 61,000, and some Windows versions use 1,025-5,000). This use indicates that

Example 3.3: Removing Unneeded Flows for Top *N*

```

<1>$ rwcut --fields=1-3,protocol,bytes --num-recs=5 traffic.rw
      sIP|                dIP|sPort|pro|        bytes|
      10.0.40.20| 192.168.70.10| 53| 17|        242|
      10.0.40.20| 192.168.70.10| 53| 17|        242|
      10.0.40.20| 192.168.70.10| 53| 17|        242|
      10.0.40.20| 192.168.70.10| 53| 17|        242|
      10.0.40.20| 192.168.70.10| 53| 17|        242|
<2>$ rwwfilter traffic.rw --aport=0,53 --fail=stdout \
| rwcut --fields=1-5,bytes --num-rec=5
      sIP|                dIP|sPort|dPort|pro|        bytes|
      10.0.40.27| 192.168.70.10|44358| 8082| 6|        332|
      10.0.40.27| 192.168.70.10|44383| 8082| 6|        332|
      10.0.40.83| 192.168.70.10|53596| 8082| 6|        838|
      10.0.40.83| 192.168.70.10|53597| 8082| 6|        551|
      10.0.40.83| 192.168.70.10|53598| 8082| 6|       1080|
<3>$ rwwfilter --start=2015/06/13 --end=2015/06/18 --type=all \
--proto=6 --dport=8082 --bytes-per=150- --pass=stdout \
| rwniq --fields=dip --values=flows,distinct:bytes
      dIP|    Records|bytes-Dist|
      155.6.3.1|         1|         1|
      192.168.200.10|       804|        32|
      192.168.70.10|      1132|        37|
<4>$ rwwfilter --start=2015/06/13 --end=2015/06/18 \
--type=all --scidr=10.0.40.83/32,192.168.200.10/32 \
--dcidr=10.0.40.83/32,192.168.200.10/32 \
--pass=traffic2.rw
<5>$ rwstats --fields=dip,dport --values=flows,bytes --count=6 \
traffic2.rw
INPUT: 11846 Records for 1954 Bins and 11846 Total Records
OUTPUT: Top 6 Bins by Records
      dIP|dPort|    Records|    Bytes| %Records|  cumul_%|
      192.168.200.10| 8082|       5906|  8252849| 49.856492| 49.856492|
      10.0.40.83|56018|         45|    6357|  0.379875| 50.236367|
      10.0.40.83|55026|         18|    4653|  0.151950| 50.388317|
      192.168.200.10| 137|         15|    3510|  0.126625| 50.514942|
      10.0.40.83| 771|         15|    2520|  0.126625| 50.641567|
      10.0.40.83|56348|          3|    3390|  0.025325| 50.666892|
<6>$ rwstats --fields=dip,dport --values=bytes,flows --count=6 \
traffic2.rw
INPUT: 11846 Records for 1954 Bins and 39548165 Total Bytes
OUTPUT: Top 6 Bins by Bytes
      dIP|dPort|    Bytes|    Records| %Bytes|  cumul_%|
      10.0.40.83|49375| 13139355|          3| 33.223678| 33.223678|
      192.168.200.10| 8082|  8252849|       5906| 20.867843| 54.091521|
      10.0.40.83|54964|  1488312|          3|  3.763290| 57.854811|
      10.0.40.83|49408|  1488312|          3|  3.763290| 61.618100|
      10.0.40.83|54345|   328470|          3|  0.830557| 62.448657|
      10.0.40.83|54404|   328470|          3|  0.830557| 63.279214|

```



192.158.200.10 is the server and 10.0.40.83 is the client. These two IP addresses account for at most a little over a third of one percent of the records. Port 137 is the Windows netbios name service port, and port 771 is an ICMP data artifact that will be discussed below.

For a contrasting look at the data, command 6 calls `rwstats` to summarize port utilization by the number of bytes (as a rough measure of the size of the communication taking place). By this measure, the largest contributor is not the traffic on port 8082, but rather traffic on an ephemeral port. This illustrates the need for analysts to examine the data from several perspectives to clarify its interpretation.

### 3.2.2 Interpreting the Top-*N* Lists

One key to interpreting the results shown in Example 3.3 is provided in the output from command 3. The last column of results is the count of distinct values for the bytes in each flow that is assigned to that bin. In this case, it is the number of bytes in each flow going to a specific IP address. For the last two entries, the value is less than 10 percent of the record count, indicating that communication with the same number of bytes is common, which in turn suggests that this is automated, rather than human-driven, traffic.

In this light, the results shown from command 5 suggest that this traffic averages about 1,400 bytes per flow to the server, with smaller acknowledgement traffic being returned to the client via the ephemeral ports. This suggestion, however, is contradicted by the results shown from command 6, which indicate that several very large-byte flows occur to the clients on the ephemeral port, indicating that data transfer is bi-directional. Confirming the data transfer dynamics and determining if any indications of threat are present requires further analysis—pulling more traffic to see if these hosts shift behavior across time as threats in their contacts to additional hosts.

In the results for command 5, the traffic to UDP port 137 (name service) is not answered with service traffic. Instead, it is responded to with messages using protocol 1, ICMPv4 (which appears with a 771 port number, although ICMP does not use ports). This is suggested by the common number of flows associated with these ports. It was confirmed by an inspection of the data using `rwcut` that was too long to show in the example. The flow generators encode the ICMP message type and code in the dPort field of NetFlow and IPFIX records.

In the example, the value of 771 corresponds to a message indicating that name service is unreachable. While the number of repetitions is not extensive (15 across 3 days of traffic), that repetition despite unreachable service indicates that the server is generating the port 137 traffic automatically.



## Chapter 4

# Intermediate Multi-path Analysis with SiLK: Explaining and Investigating

This chapter introduces intermediate multi-path analysis through application of the analytic development process with the SiLK tool suite. It discusses iteration, conditional analysis steps, categorization, and behavior identification.

Upon completion of this chapter you will be able to

- describe intermediate multi-path analysis and how it maps to the analytic development process
- describe SiLK tools commonly used with intermediate multi-path analysis
- provide example multi-path network flow analysis workflows

### 4.1 Multi-path Analysis: Concepts

#### 4.1.1 What Is Multi-path Analysis?

Some network behaviors are not visible within the single view of the network flow data provided by single-path analysis. Finding them requires investigating and integrating several different views of the data. This type of analysis is known as *multi-path analysis*.

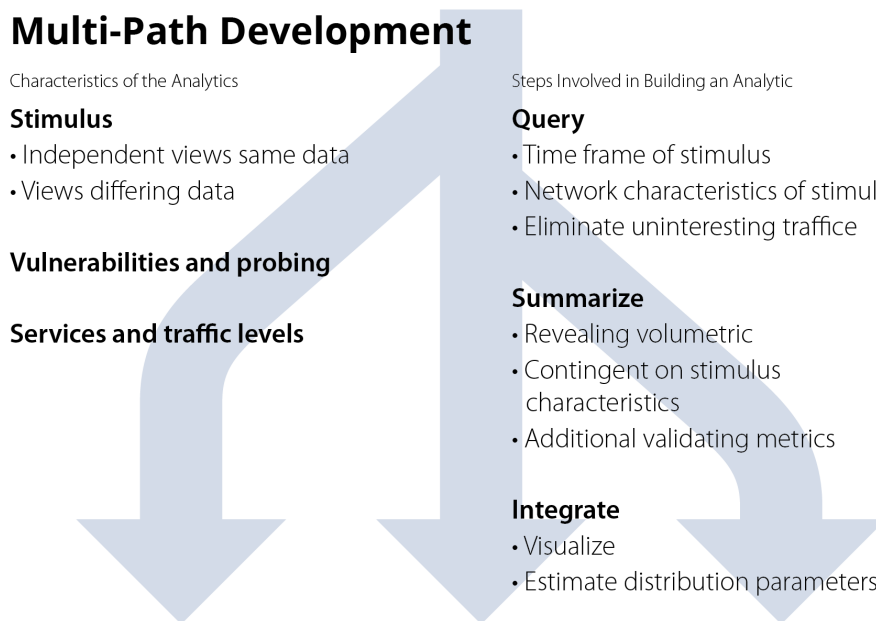
A multi-path analysis involves a deeper, multi-pronged dive into network flow data than can be accomplished with a single-path analysis. While a single-path analysis may involve looking at summary data or just one part of a data set, multi-path analysis explores different aspects of the data set (ports, IP addresses, protocols, packet and byte volumes, flow types and volumes, etc.) to find trends, leftovers, and groupings that are not necessarily visible in a single view of the data. Multi-path analysis builds upon single-path analysis; often, a single-path analysis is performed as just one phase of a multi-path analysis.

Multi-path analysis follows the general analysis framework described in Section 1.5. The overall process includes several steps:

1. Formulate the problem and gather context about the data. Is this event similar to other incidents? Which aspects of the data set should be investigated? How should we classify and group the data?
2. Model and test the data. What network behaviors are we looking for? Which statistics and metrics give us insight into these behaviors of interest?
3. Analyze the results. What did we find? How can we integrate the results of our different investigations into the data? Is our model of the event borne out by the results of our analysis?
4. Refine the analysis. Now that we have an idea of what might be going on, we can optionally take another look at the data and our assumptions to improve the analysis.

In multi-path analysis, the Formulate and Model steps are performed on multiple categories of data, each associated with aspects of the overall behavior of interest. After gathering context about the data, retrieving the data from the SiLK repository, building a model, and summarizing the data, a multi-path analysis includes an integration and analysis step that ties together the separate results to further characterize network behavior. The analyst can then iterate these steps to further refine the analysis.

Figure 4.1: Multi-Path Analysis



#### 4.1.2 Example of a Multi-path Analysis: Examining Web Service Traffic

A short example of a multi-path analysis is shown in Example 4.1. It looks at statistics on ports related to web traffic.

During the Formulate step of our multi-path analysis, we gather information about related categories of data. Multiple ports are used for web traffic: the normal HTTP port (TCP port 80), the HTTP secure

port (TCP port 443), the alternate web ports (TCP ports 8080 and 8443), and so forth. We would like to examine traffic on each of these ports.

A single-path analysis would gather flows from all of these ports into one pool of data to produce a composite model in the next step. In contrast, the multi-path analysis in Example 4.1 gathers traffic from each port as a separate pool of data. It investigates each pool individually in the Model step, profiling individual characteristics and summarizing variations.

- Command 1 uses the `mkfifo` command to create a set of named pipes, or FIFO (first-in-first-out) files, to support a complex series of `rwfilter` calls in Command 2. Each one is named after a port number (e.g., `multi-port8080.fifo`) to indicate that it carries data related to that port.

**Hint:** Named pipes efficiently transfer data from one SiLK command to another without the overhead of writing data to a disk file at each step, making the analytic run more quickly. Using named pipes also makes the analytic easier to script. Unlike regular operating system pipes (1), named pipes persist after a command finishes executing and can therefore be used as a source of input data for subsequent SiLK commands. See Appendix B.2.5 for more information.

- Commands 2 through 5 set up an analytical structure known as a *manifold*, which is discussed in more detail in section 4.2.1. The pipelined series of `rwfilter` calls in Command 2 first pulls from the repository all outbound complete TCP flows (shown in the selection and partitioning commands in the first `rwfilter` call). These outbound flows (selected by `--type=out,outweb`) include traffic produced by the monitored network. Selecting complete TCP flows (partitioned by a combination of `--proto`, `--flags-all`, `--packets`, and `--bytes-per`) avoids data that includes scans, extended sessions, redundant flow termination, and other data artifacts—all of which might confuse the analysis. The subsequent `rwfilter` calls in Command 2 then divide the retrieved flow records into separate pools, one for each of the web-related ports using both standard output and the FIFO files. Each of these `rwfilter` calls culminates (at the end of Command 2 and in Commands 3 through 5) in a call to `rwbag` to generate a summary set of byte counts per source IP address in each pool. Section 4.2.4 describes more about the bag tools.
- To make all of this work in Linux, these commands have to operate in a producer-consumer manner as background processes. Command 6 is a shell command that causes the script to wait until the producer-consumer processes have finished, so that the counts are all complete.
- Commands 7 through 10 use `rwbagcat` to calculate descriptive statistics for each of the pool of data. They send each set of statistics to a separate text file.
- Commands 11 and 12 show two of these statistics, mean and standard deviation, across the pools of data to produce integrated summaries of outbound web traffic on these ports.

### 4.1.3 Exploring Relationships and Behaviors With Multi-path Analysis

Many possible relationships in network traffic can be explored during multi-path analyses, in addition to the port-protocol alternatives in the web traffic analysis shown in Example 4.1.

- **Address relationships** can be explored by looking at the behavior of a block of addresses as seen by varying sensors, looking at variations in behavior between addresses within a block, or looking at

Example 4.1: Examining Flows for Web Service Ports

```

<1>$ mkfifo /tmp/multi-port8080.fifo; \
      mkfifo /tmp/multi-port443.fifo; \
      mkfifo /tmp/multi-port8443.fifo
<2>$ rwwfilter --start=2015/06/15 --end=2015/06/21 \
      --type=out,outweb --proto=6 --flags-all=SAF/SAF,SAR/SAR \
      --packets=4- --bytes-per=65- --pass=stdout \
| rwwfilter stdin --aport=8080 \
      --pass=/tmp/multi-port8080.fifo --fail=stdout \
| rwwfilter stdin --aport=443 --pass=/tmp/multi-port443.fifo \
      --fail=stdout \
| rwwfilter stdin --aport=8443 \
      --pass=/tmp/multi-port8443.fifo --fail=stdout \
| rwwfilter stdin --aport=80 --pass=stdout \
| rwbag --bag-file=sipv4,bytes,./output/tmp80.bag &
<3>$ rwbag --bag-file=sipv4,bytes,./output/tmp443.bag \
      /tmp/multi-port443.fifo &
<4>$ rwbag --bag-file=sipv4,bytes,./output/tmp8443.bag \
      /tmp/multi-port8443.fifo &
<5>$ rwbag --bag-file=sipv4,bytes,./output/tmp8080.bag \
      /tmp/multi-port8080.fifo &
<6>$ wait
<7>$ rwbagcat --print-stat=./output/out80.txt \
      ./output/tmp80.bag
<8>$ rwbagcat --print-stat=./output/out8080.txt \
      ./output/tmp8080.bag
<9>$ rwbagcat --print-stat=./output/out443.txt \
      ./output/tmp443.bag
<10>$ rwbagcat --print-stat=./output/out8443.txt \
      ./output/tmp8443.bag
<11>$ grep 'mean' ./output/out{80,8080,443,8443}.txt
./output/out80.txt:          mean:  5.471e+06
./output/out8080.txt:       mean:  4.707e+07
./output/out443.txt:        mean:  2.585e+07
./output/out8443.txt:       mean:  1.367e+08
<12>$ grep 'standard' ./output/out{80,8080,443,8443}.txt
./output/out80.txt:standard deviation:  2.021e+07
./output/out8080.txt:standard deviation:  6.655e+07
./output/out443.txt:standard deviation:  5.704e+07
./output/out8443.txt:standard deviation:  9.266e+06

```

behaviors of several blocks within a given Autonomous System that is handled by a common route. These address relationships might be useful for analysis tasks such as confirming suspected malware propagation or isolating targeted scanning from more general scanning.

- **Timing relationships** can be explored by separately summarizing and examining behavior before, during, and after the times associated with network events. They can also be explored by profiling behavior around an event in comparison to a period of normal activity that goes on for a similar period of time. These timing relationships could be useful to identify more subtle intrusions, for example, or to isolate activity that might indicate malicious pivoting between parts of the network.
- **Volumetric relationships** can be explored to find more complex relationships between pools of data depending on byte volumes or transfer rates. These volumetric relationships could help to indicate covert data exfiltration, for example, or detect when services are exploited to support malicious activity.

#### 4.1.4 Integrating and Interpreting the Results of Multi-path Analysis

During the Model phase of the analysis, the focus is often about determining values to provide insight on the relationships in the data. Measures of central tendency (such as averages) are frequently useful. They should be extended by measures of extent (such as range or standard deviation) to provide context for variation between pools of data. With these measures calculated, a range of activity can be specified. In the web traffic analysis shown in Example 4.1, for instance, each pool of network traffic associated with a port is profiled into an output text file using the `rwbagcat --print-stat` command, which computes a variety of descriptive statistical measures, including mean and standard deviation.

Once the measures are calculated, integrate them by matching them across the pools of data to establish trends or contrast values for the identified relationships. In Example 4.1, the standard deviations are almost all larger than the mean values, indicating that the count distribution has a long tail. More traffic values are lower than the mean than higher, but the higher ones go quite high.

Interpreting the results involves examining the statistics given and applying what insight the analyst can provide. Consider a variety of explanations for the behavior. If necessary, refine and iterate the analysis to support or deny these explanations.

Generally, consider benign interpretations for behavior first, and only reject them if appropriate contradiction can be found. In Example 4.1, the benign interpretation is that much of the traffic across these ports cannot be readily differentiated; the main difference is which port it was sent on. An alternative interpretation is that the large bulk of relatively low-byte traffic and the long tail of relatively high-byte traffic should be examined separately to determine if suspicious traffic is present. This would involve iterating the analysis to partition and separately examine low-byte and high-byte traffic for each port.

#### 4.1.5 “Gotchas” for Multi-path Analysis

While it is quite possible to explore combinations of relationships within a given analysis, some care is needed. As the number of relationships being combined increases, so does both the size and the complexity of the results. This raises several concerns in performing combined multi-path analyses:

- Interpreting the results can require a lot of rather tedious effort for limited results. It is likely that the number of data combinations that need to be evaluated will yield many cases that are not of interest to the investigation, but may still need to be explored either for completeness or to be sure

that all interesting cases are dealt with. If the combination of data relationships is not carefully chosen, analysts may spend a lot of time without much compensating insight.

- The amount of data required to fully explore combinations of relationships can be extensive, which will both slow the gathering step and require iteration across cases. Finding appropriate data to cover combinations can involve effort—for example, establishing that “normal” network activity does not itself contain malicious activity!
- As the number of combinations involved in analysis increases, so does the possibility that observed differences may occur by coincidence. This can lead to unintended “cooking” of the results, focusing on combinations that confirm the analysts’ preconceptions, rather than on a more holistic view of the behaviors.

Based on these concerns, we recommend that you **keep your multi-path analyses as simple as possible**, given the behaviors being studied. It may well be preferable to perform a series of simpler analyses rather than a single, complex, multi-factor analysis—both more manageable in action and producing more easily understood results.

## 4.2 Multi-path Analysis: Analytics

The SiLK commands, parameters, and examples described in this chapter can be employed with any analysis method. However, they involve more complex uses of the SiLK tool suite than the commands described in Chapter 2. Multi-path analyses frequently make use of these commands to construct analytics and store intermediate and final results.

### 4.2.1 Complex Filtering With `rwfilter`

Complex filtering combines operating system pipes with `rwfilter` output parameters to perform large-scale, multi-path, flow analyses. It is an important concept to adopt and incorporate into the multi-path analysis process discussed in Section 4.1 and the overall SiLK workflow described in Section 1.5.

The conditional steps, refinement, and iteration of multi-path analysis usually require multiple `rwfilter` queries. Wide time periods, large network ranges, and increasing network traffic are a few examples that complicate this process. Unfortunately, the increasing `rwfilter` directory and file searches needed to support such analyses also impact disk input/output (I/O) performance and latency.

To address these trade-offs, `rwfilter` provides three output parameters to optimize repository queries and classify traffic behavior.

- `--pass-destination` writes flow records that pass *all* partitioning criteria to a path.
- `--fail-destination` writes flow records that fail *any* partitioning criteria to an alternate path.
- `--all-destination` writes *all partitioned flow records* to a third path.

These parameters are so commonly used that the remainder of this handbook will refer to them by their common abbreviations (`--pass`, `--fail`, and `--all`). All three can be combined and repeated within the same `rwfilter` statement. They can write network flow records to a newly-created file, a named pipe, standard output (`stdout`), or standard error (`stderr`).



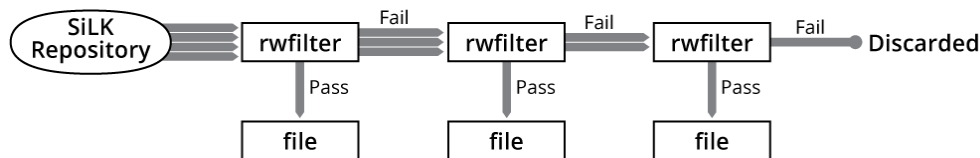
## Multi-level Filtering With Pipes and Manifolds

A *manifold* combines several `rwfilter` commands to categorize traffic. By using operating system pipes and switches with the `--pass` and `--fail` parameters, analysis can chain multiple `rwfilter` statements together to reduce an initial broad data pull into smaller sets of results that isolate traffic of interest for further analysis. After examining your manifold's initial categorization of network flow data, you can adjust the parameters of subsequent `rwfilter` calls to re-categorize data and find additional records of interest.

Manifolds perform complex traffic categorization by filtering data into overlapping and non-overlapping traffic categories.

**Manifolds with Non-overlapping Traffic.** If the categories are *non-overlapping*, the manifold uses the `--pass` parameter to write matching traffic to a file, as shown in Figure 4.2. The shaded arrows indicate data flows from the SiLK repository to the series of `rwfilter` commands that comprise the manifold. Each `rwfilter` command uses the `--pass` parameter to save the records that meet the filtering criteria to a file. It also uses the `--fail` parameter to transfer the remaining traffic to the next `rwfilter` command for further filtering. The final `rwfilter` command in Figure 4.2 saves the last category and uses the `--fail` parameter to discard the remaining uncategorized traffic. (Alternatively, the manifold could save the discarded traffic to a file.)

Figure 4.2: Diagram of a Simple, Non-overlapping Manifold



**Manifolds with Overlapping Traffic.** If the categories are *overlapping*, the manifold again uses the `--pass` parameter to filter traffic in a category. But it would use the `--all` parameter as shown in Figure 4.3 to transfer *all* traffic to the next call to `rwfilter`. With overlapping categories, the manifold is not done with a record just because it assigned a first category to it. The record may belong to other categories as well, which would be identified by subsequent calls to `rwfilter`.

Figure 4.3: Diagram of a Complex, Overlapping Manifold

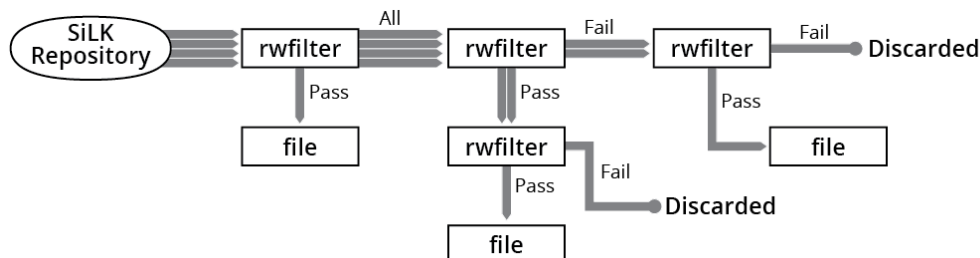


Figure 4.3 also shows how records can be passed to subsequent `rwfilter` calls via the `--fail` parameter.

Using successive `rwfilter` calls that combine the `--pass`, `--fail`, and `--all` parameters gives you a high degree of control over how your network traffic is categorized for analysis.

### Simple Manifold: Filtering Incoming Client and Server Traffic

In Example 4.2, we look for inbound flows from external servers and clients. The manifold sorts these flows into two separate files. This is an example of a non-overlapping manifold, since incoming client and server traffic have independent characteristics.

Example 4.2: Simple Manifold to Select Inbound Client and Server Flows

---

```
<1>$ rwfilter --start=2015/06/17T15 --sensor=S5 --type=in \
  --protocol=6 --flags-initial=SA/SA --packets=3- \
  --fail=stdout --pass=inbound-clients.rw \
| rwfilter stdin --flags-initial=S/SA --packets=3- \
  --pass=inbound-servers.rw
```

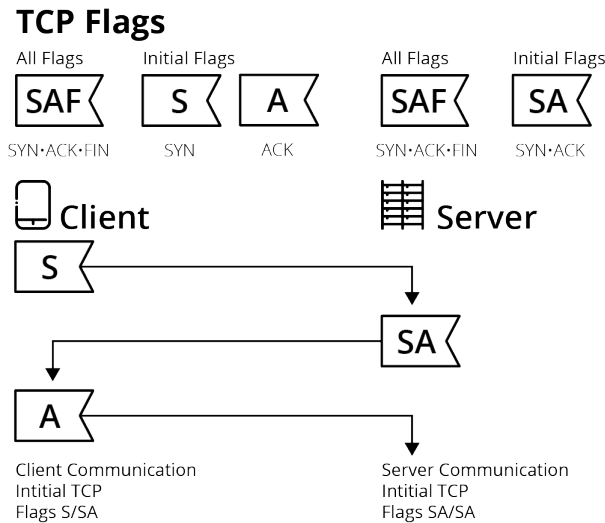
---

- The first `rwfilter` command in the manifold filters out all incoming client traffic.
  - The `--start` and `--sensor` selection parameters specify the time period and network sensor of interest.
  - `--protocol=6` selects Transmission Control Protocol (TCP) traffic.
  - `--flags-initial=SA/SA` selects traffic where TCP flags on the first packet of the flow record are set to “SYN/ACK”, indicating that the source IP address is a server.
  - `--packets=3-` selects flow records with three or more packets to identify TCP sockets.
  - `--pass=incoming-client.rw` stores flow records that match the filtering criteria to the file `incoming-client.rw`.
  - `--fail=stdout` sends all records that do not match the filtering criteria to the second part of the manifold.

**Hint:** Figure 4.4 shows the TCP flags for client and server communication. The client first sends a SYN flag to initiate the TCP connection. The server responds with a SYN flag to initiate its own TCP connection with the client and an ACK flag to acknowledge the client’s request for a TCP connection. For more detailed information about TCP flags, see Appendix A.5.2.

- The second `rwfilter` command in the manifold selects incoming server traffic.
  - `stdin` tells the `rwfilter` command to process flow records from standard input (i.e., the output from the previous `rwfilter` command).
  - The command does not filter on Sensor, Type and Protocol. It looks at traffic that does not meet those criteria.
  - `--flags-initial=S/SA` selects traffic where TCP flags on the first packet of the flow record are set to “SYN” with no “ACK,” indicating the source IP address is a client.
  - `--packets=3-` selects flow records with three or more packets, which is the minimum number required to open a TCP socket.

Figure 4.4: Client and Server TCP flags



- `--pass=incoming-server.rw` stores flow records that match the filtering criteria to the file `incoming-server.rw`
- Notice that the second `rwfilter` command did not specify a `--fail` destination. We are only interested in the records that pass the filtering criteria in the second part of the manifold, so there is no need to save the ones that fail.

### Expanding the Simple Manifold: Filtering for Incoming and Outgoing Client and Server Traffic

To expand the simple manifold in the previous example to partition both inbound and outbound traffic, we make three modifications as shown in Example 4.3.

1. Prefilter traffic to ignore flows that we know are unwanted. These consist of non-TCP flow types other than `in` and `out` plus flows with fewer than three packets.
  2. Filter traffic from internal and outgoing servers and clients into files.
  3. Store the leftover flows in a file for later use.
- The first `rwfilter` command in the manifold filters out unwanted traffic.
    - `--start` and `--sensor` selection specify the time period and network sensor of interest.
    - `--protocol=6` selects Transmission Control Protocol (TCP) traffic.
    - `--packets=3-` selects flow records with three or more packets indicating TCP sockets.
    - `--pass=stdout` sends all records that match the filtering criteria to the second part of the manifold. Records that do not match are ignored.
  - The second `rwfilter` command in the manifold filters out inbound server traffic.

Example 4.3: Complex Manifold to Select Inbound Client and Server Flows

```

<1>$ rffilter --start=2015/06/17T15 --sensor=S5 --type=in,out \
  --protocol=6 --packets=3- --pass=stdout \
| rffilter stdin --type=in --flags-initial=SA/SA \
  --pass=incoming-server.raw --fail=stdout \
| rffilter stdin --type=in --flags-initial=S/SA \
  --pass=incoming-client.raw --fail=stdout \
| rffilter stdin --type=out --flags-initial=SA/SA \
  --pass=outgoing-server.raw --fail=stdout \
| rffilter stdin --type=out --flags-initial=S/SA \
  --pass=outgoing-client.raw --fail=leftover.raw
<2>$ for f in incoming-server.raw incoming-client.raw \
  outgoing-server.raw outgoing-client.raw leftover.raw; \
  do echo -n "$f: "; \
  rffileinfo --fields=count-records $f; \
  done
incoming-server.raw: incoming-server.raw:
  count-records      1001
incoming-client.raw: incoming-client.raw:
  count-records      41
outgoing-server.raw: outgoing-server.raw:
  count-records      41
outgoing-client.raw: outgoing-client.raw:
  count-records      1002
leftover.raw: leftover.raw:
  count-records      2

```

`stdin` tells the `rwfilter` command to process flow records from standard input (i.e., the output from the previous `rwfilter` command).

`--type=in` selects inbound traffic.

`--flags-initial=SA/SA` selects traffic where TCP flags on the first packet of the flow record are set to "SYN/ACK", indicating that the source IP address is a server.

`--pass=incoming-server.raw` stores flow records that match the filtering criteria to the file `incoming-server.raw`

`--fail=stdout` sends all records that do not match the filtering criteria to the next part of the manifold.

- The third `rwfilter` command in the manifold filters out inbound client traffic.

`--type=in` selects inbound traffic.

`--flags-initial=S/SA` selects traffic where TCP flags on the first packet of the flow record are set to "SYN", indicating the source IP address is a client.

`--pass=outgoing-client.raw` stores flow records that match the filtering criteria to the file `outgoing-client.raw`.

- The fourth `rwfilter` command in the manifold filters out outbound server traffic.

`--type=out` selects outbound traffic.

`--flags-initial=SA/SA` selects traffic where TCP flags on the first packet of the flow record are set to "SYN/ACK", indicating the source IP address is a server.

`--pass=outgoing-server.raw` stores flow records that match the filtering criteria to the file `outgoing-server.raw`

- The fifth (and final) `rwfilter` command in the manifold filters out outbound client traffic.

`--type=out` selects outbound traffic.

`--flags-initial=S/SA` selects traffic where TCP flags on the first packet of the flow record are set to "SYN", indicating the source IP address is a client.

`--pass=outgoing-client.raw` stores flow records that match the filtering criteria to the file `outgoing-client.raw`

`--fail=leftover.raw` saves all records that do not match the filtering criteria in the file `leftover.raw`.

## 4.2.2 Finding Low-Packet Flows with `rwfilter`

The TCP state machine is complex (see Figure A.4). As described in Appendix A.5, legitimate service requests require a minimum of three (more commonly four) packets in the flow from client to server. Flows from server to client may only have two packets.

Several types of illegitimate traffic (such as port scans and responses to spoofed-address packets) involve TCP flow records with low numbers of packets. Although legitimate TCP flow records occasionally have low numbers of packets (such as continuations of previously timed-out flow records, contact attempts on hosts that do not exist or are down, services that are not configured, and RST packets on already closed connections), this behavior is relatively rare. Understanding where low-packet TCP network traffic comes

from and when such flow records are collected most frequently can therefore help you to identify traffic that is potentially illegitimate.

Example 4.4 shows how to use the manifold concept from Section 4.2.1 to find low-packet traffic. It illustrates how `rwfilter` can be used to refine selections to isolate flow records of interest.

1. The first call to `rwfilter` in command 1 selects all incoming flow records in the repository on 6/17/2015, that describe TCP traffic, and that had one to three packets in the flow record. It is the only `rwfilter` call that pulls data directly from the SiLK repository; as such, it is the only one that uses selection parameters.

This call also uses a combination of partitioning parameters (`--protocol` and `--packets`) to isolate low-packet TCP flow records from the selected time range. It uses the `--pass` switch twice: once to save the selected records to the file `./Ex4-data/lowpacket.rw` and once to direct the selected records to standard output (`stdout`) for use by the next `rwfilter` command in the manifold.

The `--print-statistics` parameter saves information about the `rwfilter` call to the file `temp-all.txt`, including the number of repository files retrieved by `rwfilter`, the total number of flow records, and the number of records that passed or failed the filter.

2. The second call to `rwfilter` in command 1 uses `--flags-all` as a partitioning parameter to pull out flow records of interest. It passes flow records that had the SYN flag set in any of their packets, but do not have the ACK, RST, and FIN flags set in any of their packets. It fails those that did not show this flag combination. It saves statistics to the file `temp-syn.txt`. Records that fail this filter are passed to the next `rwfilter` call via the `--fail=stdout` parameter.
3. The third call to `rwfilter` extracts the flow records that have the RST flag set, but had the SYN and FIN flags not set. It saves statistics to the file `temp-rst.txt`.
4. Command 2 displays the statistics information saved at each step in the manifold. We can see how the succession of calls to `rwfilter` progressively refine the data, and how data passes from one filter to the next.

#### Example 4.4: Extracting Low-Packet Flow Records

---

```
<1>$ rwfilter --start-date=2015/06/17 \
  --type=in,inweb --protocol=6 --packets=1-3 \
  --print-statistics=temp-all.txt \
  --pass=./Ex4-data/lowpacket.rw --pass=stdout \
| rwfilter --flags-all=S/SARF \
  --print-statistics=temp-syn.txt \
  --pass=./Ex4-data/synonly.rw --fail=stdout stdin \
| rwfilter --flags-all=R/SRF \
  --print-statistics=temp-rst.txt \
  --pass=./Ex4-data/reset.rw stdin
<2>$ cat temp-all.txt output/temp-syn.txt output/temp-rst.txt
Files   415.  Read   9083613.  Pass   906016.  Fail   8177597.
Files    1.  Read   906016.  Pass    97279.  Fail   808737.
Files    1.  Read   808737.  Pass   130145.  Fail   678592.
```

---

### 4.2.3 Time Binning, Options, and Thresholds With `rwstats`, `rwuniq` and `rwcount`

#### Approximating Flow Behavior Over Time

SiLK flow records do not contain information about the time distribution of packets and bytes during a flow. Grouping packets into flow records results in a loss of timing information; specifically, it is not possible to tell how the packets in a flow are distributed over time. Even at the sensor, the information about the time distribution of packets in a flow is lost.

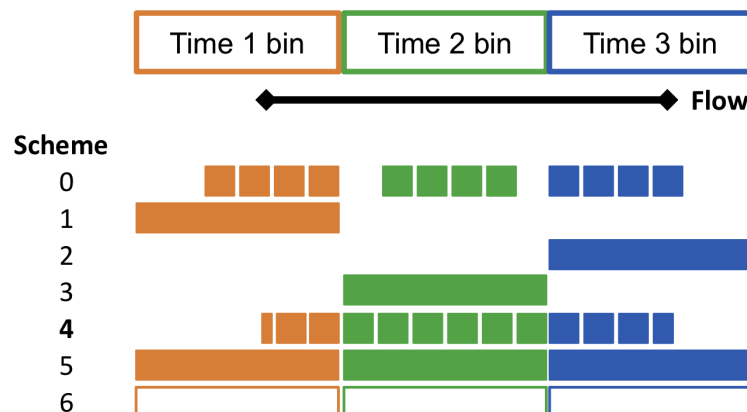
By default, SiLK distributes the packets and bytes equally across all the milliseconds in the flow's duration. This approximation works well for investigating overall trends in network behavior.

However, some types of analysis benefit from intentionally changing the time distribution of packets. For example, incident analysis investigates behavior during specific time bands. Changing the time distribution of packets and bytes (or *load-scheme*) emphasizes different flow characteristics of interest.

Analysts can impose a specific time distribution on `rwcount` by using the `--load-scheme` parameter. `rwcount` can assign one of seven time distributions of packets and bytes in the flow to allocate the volume to time bins.

Figure 4.5 illustrates the allocation of flows, packets, and bytes to time bins under different load-schemes (which are described more fully in Table C.12). The squares depict a fixed number of packets or bytes. The partial squares are proportional to the complete squares. The wide, solidly filled rectangles depict entire flows, along with their packets and bytes. They appear once in schemes 1, 2, and 3 (where one bin receives the entire flow with its packets and bytes) and in every bin for scheme 5 (where every bin receives the entire flow). The wide, hollow rectangles only appear in scheme 6 and represent whole flows with no packets or bytes.

Figure 4.5: Allocating Flows, Packets and Bytes via `rwcount` Load-Schemes



## Using Thresholds to Profile a Slice of Flows

`rwuniq` allows you to set thresholds to segregate IP addresses by the number of flows, sizes of flows, and other values. Recall that `rwuniq` reads SiLK flow records, groups them according to a key composed of user-specified flow attributes, then computes summary values for each group (or bin), such as the sum of the bytes fields for all records that match the key. Thresholding limits the output of `rwuniq` to bins where the summary value meets user-specified minimum or maximum values.

The `--bytes`, `--packets`, and `--flows` parameters are all threshold operators for filtering. For example, to show only the source IP addresses with 200 flow records or more, use the `--flows=200-` parameter as shown in Example 4.5.

Example 4.5: Constraining Counts to a Threshold by using `rwuniq --flows`

```
<1>$ rfilter --start-date=2015/06/17T15 --type=in --protocol=0- \
--pass=./Ex4-data/in_month.rw
<2>$ ls -l ./Ex4-data/in_month.rw
-rw-r--r--. 1 analyst analyst 7409538 Apr 25 17:22 ./Ex4-data/in_month.rw
<3>$ rwuniq ./Ex4-data/in_month.rw --field=sIP --value=Flows \
--flows=200- \
| head -n 10
      sIP |   Records |
192.168.143.57 |     481 |
192.168.161.26 |     443 |
  192.168.40.51 |     254 |
    10.0.40.23 |    5128 |
192.168.165.83 |     378 |
  192.168.40.25 |   11951 |
192.168.162.160 |     401 |
    10.0.40.92 |    1838 |
192.168.143.162 |     635 |
```

In addition, `rwuniq` can count bytes and packets for a flow threshold through the `bytes` and `packets` values in the `--values` parameter, as shown in Example 4.6. This example counts the byte and packet volumes for all IP addresses that exceed a minimum flow threshold of 2000 records (`--values=Bytes, Packets, Flows --flows=2000-`).

If a range (such as `--flows=2000-`) is not specified, the parameter simply adds the named count to the list started by the `--values` parameter. We recommend using the `--values` parameter for this purpose. `--values` provides greater control over the order in which the values are displayed than the other thresholding parameters.

**Hint:** `rwuniq` provides three ways to specify ranges: with the low and high bounds separated by a hyphen (e.g., 200–2000); with a low bound followed by a hyphen (e.g., 200–) denoting that there is no upper bound; and with a low bound alone (e.g., 200). Unlike `rfilter` partitioning values, the last method denotes a range with no upper bound, not just a single value. **We do not recommend using this method because it can lead to confusion.**

If multiple threshold parameters are specified, `rwuniq` will print all records that meet *all* of the threshold criteria, as shown in Example 4.7.



Example 4.6: Setting Minimum Flow Thresholds with `rwuniq --values`

```

<1>$ rfilter --start-date=2015/06/17T15 --type=in --protocol=0- \
  --pass=./Ex4-data/in_month.rw
<2>$ ls -l ./Ex4-data/in_month.rw
-rw-r--r--. 1 analyst analyst 7409538 Apr 25 17:22 ./Ex4-data/in_month.rw
<3>$ rwuniq ./Ex4-data/in_month.rw --field=sIP \
  --values=Bytes, Packets, Flows --flows=2000- \
| head -n 10

```

sIP	Bytes	Packets	Records
10.0.40.23	9331109	108622	5128
192.168.40.25	6066944	58819	11951
192.168.20.58	4128957	54314	28189
10.0.40.53	64029248	281180	33917
192.168.200.10	6089612	20535	8816
10.0.40.54	6131060	42785	7075
10.0.40.20	36120528	345466	149284
10.0.20.58	4043680	55644	21253
67.215.0.8	1499433280	6096359	30549

Example 4.7: Constraining Flow and Packet Counts with `rwuniq --flows` and `--packets`

```

<1>$ rfilter --start-date=2015/06/17T15 --type=in --protocol=0- \
  --pass=./Ex4-data/in_month.rw
<2>$ ls -l ./Ex4-data/in_month.rw
-rw-r--r--. 1 analyst analyst 7409538 Apr 25 17:23 ./Ex4-data/in_month.rw
<3>$ rwuniq ./Ex4-data/in_month.rw --field=sIP \
  --values=Bytes, Packets, Flows --flows=2000- \
  --packets=100000-

```

sIP	Bytes	Packets	Records
10.0.40.23	9331109	108622	5128
10.0.40.53	64029248	281180	33917
10.0.40.20	36120528	345466	149284
67.215.0.8	1499433280	6096359	30549
192.168.40.20	67729921	283002	106848

## Profiling Using Compound Keys

Profiling can also be done by counting and thresholding combinations of fields, in addition to the simple counting shown previously. To use a compound key, specify it using a comma-separated list of values or ranges in the `rwuniq --fields` parameter. Keys can be manipulated in the same way as with `rwcut`: `--fields=3,1` is a different key than `--fields=1,3`.

In Example 4.8, the `--fields` parameter is used to identify communication between clients and specific services only when the number of flows for the key exceeds a threshold. It counts and thresholds incoming traffic to identify those source IP addresses with the highest number of flow records that connect to specific TCP ports (`--fields=sIP,sPort`).

Alternatively, you can use the `rwstats` command with compound keys to perform this analysis. In Example 4.9, the `--fields` parameter is similarly used to count incoming traffic and identify the eleven source IP

Example 4.8: Profiling IP addresses with `rwuniq --fields`

```

<1>$ rfilter --start-date=2015/06/17T15 --type=in --protocol=6 \
--pass=./Ex4-data/in_month.rw
<2>$ ls -l ./Ex4-data/in_month.rw
-rw-r--r--. 1 gtsanders domain users 2059092 May 14 12:15 ./Ex4-data/in_month.rw
<3>$ rwuniq --fields=sIP,sPort --value=Flows --flows=20- \
./Ex4-data/in_month.rw \
| head -n 11

```

sIP	sPort	Records
192.168.143.162	591	26
192.168.111.131	591	28
192.168.122.141	591	26
10.0.40.20	88	476
10.0.40.20	139	1189
192.168.164.119	591	26
192.168.40.20	60309	20
192.168.40.25	445	120
192.168.165.216	591	26
192.168.40.20	135	265

addresses with the highest number of flow records that connect to specific TCP ports (`--fields=sIP,sPort`). In addition to counting and displaying these records, `rwstats` computes their cumulative statistics. This allows you to directly compare the amount of traffic carried by each source IP-port combination.

### Isolating Behaviors of Interest

`rwuniq` can be used in conjunction with `rfilter` to profile flow records for a variety of behaviors:

1. Use `rfilter` to filter records for the behavior of interest. (To filter for multiple behaviors, set up a manifold as described in Section 4.2.1.)
2. Use `rwuniq` to count the records that exhibit that behavior.

This can help you to understand the behavior of hosts that use or provide a variety of services. Example 4.10 shows how to generate data that compare hosts showing DNS (domain name system) and non-DNS behavior among a group of flow records. We can find DNS servers by filtering the file `in_month.rw` (created in Example 4.8) for hosts that carry traffic on port 53 with the UDP protocol (17).

1. Command 1 first isolates the set of hosts of interest by using `rfilter` to filter records from `in_month.rw` with UDP traffic on port 53 (`--protocol=17 --aport=53`). It then uses `rwset` to generate an IPset (`interest.set`) from the records that pass the filter.
2. Command 2 uses `interest.set` to filter `in_month.rw` again to distinguish IP addresses with general UDP traffic from the set of hosts that carry DNS traffic on port 53. It then uses `rwuniq` to count the DNS flow records and sorts them by source address.

Although `rwuniq` will correctly sort the output rows by IP address without zero-padding, the upcoming `join` command will not understand that the input is properly sorted without `rwuniq` first preparing the addresses with the `--ip-format=zero-padded` parameter.

Example 4.9: Profiling IP addresses with `rwstats --fields`

```

<1>$ rwwfilter --start-date=2015/06/17T15 --type=in --protocol=6 \
  --pass=./Ex4-data/in_month.rw
<2>$ ls -l ./Ex4-data/in_month.rw
-rw-r--r--. 1 analyst analyst 2059092 May 14 12:26 ./Ex4-data/in_month.rw
<3>$ rwstats --count=11 --fields=sIP,sPort --value=Flows \
  ./Ex4-data/in_month.rw \
INPUT: 124356 Records for 11637 Bins and 124356 Total Records
OUTPUT: Top 11 Bins by Records

```

sIP sPort	Records	%Records	cumul_%
10.0.40.53  5723	33103	26.619544	26.619544
67.215.0.8 11009	12479	10.034900	36.654444
67.215.0.8 11007	10714	8.615588	45.270031
67.215.0.8  135	7181	5.774550	51.044582
10.0.40.23  8443	5128	4.123645	55.168227
192.168.40.20  88	3296	2.650455	57.818682
10.0.40.20  445	1986	1.597028	59.415710
10.0.40.20  389	1367	1.099263	60.514973
10.0.40.20  139	1189	0.956126	61.471099
192.168.70.10  8082	1077	0.866062	62.337161
10.0.40.20 49158	1071	0.861237	63.198398

3. Command 3 counts the non-DNS flow records created with the `--fail` switch in Command 2 and sorts them by source address.
4. Command 4 merges the two count files by source address and then sorts them by number of DNS flows with the results shown. Hosts with high counts in both columns should be either workstations or gateways. Hosts with high counts in DNS and low counts in non-DNS should be DNS servers.<sup>6</sup>

For more complex summaries of behavior, use the `rwbag` command and its related utilities as described in Section 4.2.4.

#### 4.2.4 Summarizing Network Traffic with Bags

IPsets contain lists of IP addresses. However, it's often useful to associate a value with each address in an IPset. For instance, you may want to associate the IP addresses that engage in web traffic with the volume of flows, packets, or bytes of web traffic that each address carries. *Bags* are extended sets that contain these types of key-value pairs.

Where IPsets record the presence or absence of key values, bags add the ability to count the number of instances of a particular key value—that is, the number of bytes, the number of packets, or the number of flow records associated with that key. Bags also allow the analyst to summarize traffic on characteristics other than IP addresses—specifically on protocols and ports.<sup>7</sup>

Bags can be thought of as enhanced IPsets. Like IPsets, they are binary structures that can be manipulated using a collection of tools. As a result, operations that are performed on IPsets have analogous bag operations,

<sup>6</sup>The full analysis to identify DNS servers is more complex and will not be dealt with in this handbook.

<sup>7</sup>PySiLK allows for even more general bag key values and count values. See the documentation *PySiLK: SiLK in Python* for more information.

Example 4.10: Isolating DNS and Non-DNS Behavior with `rwuniq`

```

<1>$ rfilter in_month.rw --protocol=17 --aport=53 --pass=stdout \
| rset --sip-file=interest.set
<2>$ rfilter in_month.rw --sipset=interest.set --protocol=17 \
--pass=stdout \
| rfilter stdin --aport=53 --fail=not-dns.rw --pass=stdout \
| runiq --fields=sIP --no-titles --ip-format=zero-padded \
--sort-output --output-path=dns-saddr.txt
<3>$ runiq not-dns.rw --fields=sIP --no-titles \
--ip-format=zero-padded --sort-output \
--output-path=not-dns-saddr.txt
<4>$ echo '          sIP|          DNS||   not DNS||' \
; join -t'|' dns-saddr.txt not-dns-saddr.txt \
| sort -t'|' -nrk2,2 \
| head -n 5
          sIP|          DNS||   not DNS|
010.000.040.020|    124652||    14322|
192.168.040.020|    98128||     797|
192.168.200.010|    7123||     64|
192.168.040.025|    5188||    5127|
192.168.165.216|     508||     47|

```

such as addition (the equivalent to union). Analysts can also extract a cover set (the set of all IP addresses in the bag) for use with `rfilter` and the IPset tools.

### Generating Bags from Network Flow Data

The `rbag` command creates bags from raw network flow data, either directly output from the `rfilter` command or stored in a file. The `key` parameter specifies the network flow record field that serves as the key for the bag. Examples of keys include source IP address (`sIPv4`, `sIPv6`), destination IP address (`dIPv4`, `dIPv6`), source port (`sPort`), destination port (`dport`), `protocol`, `packets`, and `bytes`.

The `counter` parameter sums up the number of `records`, `flows`, `packets`, or `bytes` for the flow record field specified by `key`. `outputfile` is the name of the file where the bag is stored, such as `mybag.bag`. Bags are stored in binary format to make analysis tasks faster and more efficient. Use the `rbagcat` command to display the contents of a bag.

Example 4.11 shows an example of how to use `rbag` in conjunction with `rfilter`. You may specify multiple `--bag-file` parameters when you issue the `rbag` command.

Example 4.11: Generating Bags with `rbag`

```

<1>$ rfilter --type=in,inweb --start-date=2015/06/18T12 \
--protocol=6 --pass=stdout \
| rbag --sip-packets=x.bag --dip-flows=y.bag
<2>$ file x.bag y.bag
x.bag: data
y.bag: data

```

See Appendix C.18 for more information about the `rwbag` command. To view a complete list of command options, type `rwbag --help` at the command line.

### Summarizing Network Traffic with Bags

To show how useful bags can be, we return to the task mentioned earlier in this section: analyzing outgoing web traffic. We want to find out which IP addresses engage in web traffic and will narrow our study to outgoing TCP flows on ports 80 and 443. We can find these IP addresses by using the `rwfilter` and `rwuniq` commands as shown in Example 4.12.

Example 4.12: Summarizing Network Traffic with `rwuniq`

```
<1>$ rwfilter --start-date=2015/06/17 --sensors=S1 --type=outweb \
  --protocol=6 --sport=80,443 --packets=3- --pass=stdout \
| rwuniq --fields=sip --values=bytes
      sIP|          Bytes|
192.168.40.24|      877782|
192.168.20.59|     1392516|
192.168.40.91|     124548|
192.168.40.92|     124548|
```

This provides us with addresses and byte counts in text format. However, we would like to use this information during further analysis with SiLK—for example, as an IPset with the addresses and some way to store the number of bytes for each address. To store such a list, we need to create a bag with the `rwbag` command as shown in Example 4.13. We want the key to be the IP address and the value to be the number of bytes of outbound web traffic.

Example 4.13: Summarizing Network Traffic with Bags

```
<1>$ rwfilter --start-date=2015/06/17 --sensors=S1 --type=outweb \
  --protocol=6 --sport=80,443 --packets=3- --pass=stdout \
| rwbag --bag-file=sipv4,sum-bytes,outgoingweb.bag \
<2>$ file outgoingweb.bag
outgoingweb.bag: SiLK, RWBAG v3, Little Endian, LZ0 compression
<3>$ rwbagcat outgoingweb.bag
192.168.20.59|      1392516|
192.168.40.24|      877782|
192.168.40.91|     124548|
192.168.40.92|     124548|
```

The file `outgoingweb.bag` contains the list of IP addresses that carry outgoing web traffic and the volume of outbound traffic in bytes that flows through each address. Unlike the output of the `rwuniq` command, this information is stored in a single binary file, as shown in 4.13. We can now use this file during further analysis of outbound web traffic.

### Generating Bags From IP Sets or Text: A Scanning Example

You can create a bag from an existing set or a text file by using the `rwbagbuild` tool. This allows you to associate counts with items in a set or file. It also gives you more flexibility with creating bags than the

`rwbag` command does. For instance, you can use `rwbagbuild` to count something other than bytes, packets, or flow records for an address.

`rwbagbuild` takes either an IPset (as specified in `--set-input`) or a text file (as specified in `--bag-input`), but not both.

- For IPset input, the `--default-count` parameter specifies the count value for each set element in the output bag. If no `--default-count` value is provided, the count will be set to one.
- For text-file input, the lines of the file are expected to consist of a key value, a delimiter (by default the vertical bar), and a count value. Keys can be IP addresses (including canonical forms, CIDR blocks, and SiLK address wildcards) or unsigned integers.

See Appendix C.19 for more information about the `rwbagbuild` command. To view a complete list of command options, type `rwbagbuild --help` at the command line.

Example 4.14 shows how to use the `rwbagbuild` command in conjunction with the `rwscan` command to create a bag that contains IP addresses that show evidence of scanning activity and the number of flows associated with them.

`rwscan` analyzes SiLK flow records for signs of network scanning—when an external host gathers information about a network during the reconnaissance phase of an attack. It takes sorted network flow records as input and outputs in columnar text format any IP addresses that show signs of network scanning. This is useful for identifying hosts that are conducting reconnaissance and the ports and protocols of interest to them. Pairing this command with `rwbagbuild` allows you to create a bag that stores scanner IP addresses and attributes of their activity for further investigation.

1. Command 1 uses `rwfilter` to pull inbound TCP traffic (`--proto=6 --type=in,inweb`).
2. The `rwscan` command requires input to be pre-sorted by source IP address, protocol, and destination IP address. Command 1 therefore calls `rwsort --fields=sip,proto,dip` to sort the selected records.
3. Command 1 then uses `rwscan` to search for IP addresses that show signs of scanning activity. It pipes the output through the operating system `cut` command to remove the delimiters (`|`) in the `rwscan` output.
4. Finally, command 1 use the `rwbagbuild` command to create a bag (`scanners.bag`) from the `rwscan` output. It uses the scanning IP addresses as the key and the number of flow records as the associated count for the bag entries.
5. Commands 2 and 3 display the list of scanners created in command 1. Command 2 uses the `rwbagcat` command to create a text file that contains the contents of `scanners.bag`. (`rwbagcat` is described later in this section.) Command 3 shows that file.

**Specifying IP addresses with `rwbagbuild`.** The `rwbagbuild` command does not support mixed input of IP addresses and integer values, since there is no way to specify whether the number represents an IPv4 address or an IPv6 address. (For example, does 1 represent `::FFFF.0.0.0.1` or `::1`?) `rwbagbuild` also does not support symbol values in its input, so some types commonly expressed as symbols (TCP flags, attributes) must be translated into an integer form.

Example 4.14: Creating a Bag of Network Scanners with `rwbagbuild` and `rwscan`

```

<1>$ rwfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --proto=6 --type=in,inweb --pass=stdout \
| rwsort --fields=sip,proto,dip \
| rwscan --scan-model=2 --output-path=stdout --no-title \
| cut -f1,5 -d'|' \
| rwbagbuild --bag-input=stdin --key-type=sIPv4 \
  --counter-type=records >scanners.bag
<2>$ echo 'Scanner | Flows|' >scanners.txt \
; rwbagcat scanners.bag >>scanners.txt
<3>$ cat scanners.txt
Scanner | Flows|
192.168.181.8| 45428|

```

Similarly, `rwbagbuild` does not support formatted time strings. Times must be expressed as unsigned integer seconds since UNIX epoch. If the delimiter character is present in the input data, it must be followed by a count. If the `--default-count` parameter is used, its argument will override any counts in the text-file input; otherwise the value in the file will be used. If no delimiter is present, either the `--default-count` value will be used or the count will be set to 1 if no such parameter is present. If the key value cannot be parsed or a line contains a delimiter but no count, `rwbagbuild` prints an error and exits.

### Displaying the Contents of Bags

To view or summarize the contents of a bag, use the `rwbagcat` command. By default, it displays the contents of a bag in sorted order as shown in Example 4.15.

Example 4.15: Viewing the Contents of a Bag with `rwbagcat`

```

<1>$ rwbagcat x.bag \
  | head -n 5
192.0.2.198| 1281|
192.0.2.227| 12|
192.0.2.249| 90|
198.51.100.227| 3|
198.51.100.244| 101|

```

See Appendix C.20 for more information about the `rwbagcat` command. To view a complete list of command options, type `rwbagcat --help` at the command line.

**Thresholding Bags.** In Example 4.15, the *counts* (the number of elements that match a particular IP address) are printed per key. `rwbagcat` can also print values within ranges of both counts and keys, as shown in Example 4.16.

These thresholding values can be used in any combination.

Example 4.16: Thresholding Results with `rwbagcat --mincounter, --maxcounter, --minkey, and --maxkey`

```
<1>$ rwbagcat --mincounter=100 --maxcounter=600 kerbserv.bag
 10.0.40.20|          574|
 67.215.0.5|         245|
<2>$ rwbagcat --minkey=10.0.0.0 --maxkey=192.168.255.255 \
 kerbserv.bag
 10.0.40.20|          574|
 67.215.0.5|         245|
192.168.40.20|        4596|
```

**Counting Keys in Bags.** In addition to thresholding, `rwbagcat` can also reverse the index; that is, instead of printing the number of counted elements per key, it can produce a count of the number of keys matching each count using the `--bin-ips` parameter. This reverse count is shown in Example 4.17.

- In Command 1, it is shown using linear binning—one bin per value, with the counts showing how many keys had that value.
- In Command 2, it is shown with binary binning—values collected by powers of two and with counts of keys having summary volume values in those ranges.
- In Command 3, it is shown with decimal logarithmic binning—values collected in bins that provide one bin per value below 100 and an even number of bins for each power of 10 above 100, arranged logarithmically and displayed by midpoint. This option supports logarithmic graphing options.

Example 4.17: Displaying Unique IP Addresses per Value with `rwbagcat --bin-ips`

```
<1>$ rwbagcat --bin-ips dns.bag \
| head -n 5
          1|          1|
          3|          1|
         14|          1|
         18|          1|
         30|          1|
<2>$ rwbagcat --bin-ips=binary dns.bag \
| head -n 5
 2^00 to 2^01-1|          1|
 2^01 to 2^02-1|          1|
 2^03 to 2^04-1|          1|
 2^04 to 2^05-1|          2|
 2^09 to 2^10-1|          3|
<3>$ rwbagcat --bin-ips=decimal dns.bag \
| head -n 5
          1|          1|
          3|          1|
         14|          1|
         18|          1|
         30|          1|
```



The `--bin-ips` parameter can be particularly useful for distinguishing between sites that are hit by scans (where only one or two packets may appear) from sites that are engaged in legitimate activity.

**Formatting Key Values for Bags.** If the bag is not keyed by IP address, the optional `--key-format` parameter makes it *much* easier to read the output of `rwbagcat`. Example 4.18 shows the difference in output for a sIP-keyed bag counting bytes, where the IP addresses are shown in decimal and hexadecimal formats.

Example 4.18: Displaying Decimal and Hexadecimal Output with `rwbagcat --key-format`

---

```
<1>$ rwbagcat kerbserv.bag
  10.0.40.20|          751382|
  67.215.0.5|         395218|
 192.168.40.20|       5510424|
<2>$ rwbagcat --key-format=decimal kerbserv.bag
167782420|          751382|
1138163717|         395218|
3232245780|       5510424|
<3>$ rwbagcat --key-format=hexadecimal kerbserv.bag
 a002814|          751382|
43d70005|         395218|
c0a82814|       5510424|
```

---

### Comparing the Contents of Bags

Once you have created bags to store key-value pairs, you can compare their contents to identify common values and trends. For example, you may want to compare the traffic volumes associated with two groups of IP addresses, each in a separate bag file. Use the `rwbagtool --compare` parameter to compare the contents of two bags. It stores the output from this comparison in a new bag file.

For each *key* that appears in both bag files, the `--compare` option compares the value of the key's associated *counter* (i.e., the number of bytes, packets, or records summed up by the `rwbag` or `rwbagbuild` command) in the first file to the value of the key's counter in the second file.

- If the comparison is *true*, the key appears in the resulting bag file with a counter of 1.
- If the comparison is *false*, the key is not present in the output file.
- Keys that appear in only one of the input bag files are ignored.

`rwbagtool --compare` can perform the following comparisons:

- `lt` Finds keys in the first bag file whose counters are less than those in the second bag file.
- `le` Finds keys in the first bag file whose counters are less than or equal to those in the second bag file.
- `eq` Finds keys whose counters are equal in both files.
- `ge` Finds keys in the first bag file whose counters are greater than or equal to those in the second bag file.

`gt` Finds keys in the first bag file whose counters are greater than those in the second bag file.

See Appendix C.21 for more information about the `rwbagtool` command. To view a complete list of command options, type `rwbagtool --help` at the command line.

## 4.2.5 Working with Bags and IPsets

Since bags are essentially enhanced IPsets, SiLK provides operations that enable you to create IPsets from bags and compare the contents of bags with those of IPsets.

### Extracting IPsets from Bags

Sometimes you will want to extract an IPset from a bag—for instance, if you used the `rwbagtool --compare` command to compare traffic associated with IP addresses in two bags and want to analyze the set of IP addresses in the new bag. Use the `rwbagtool --coverset` parameter to generate a *cover set*: the set of IP addresses in a bag. The resulting IPset file can be used with `rwfilter` and manipulated with any of the `rwset` commands.

Example 4.19 shows how to extract an IPset from a bag file. The `rwsetcat` command displays the contents of the resulting set and the `rwbagcat` command displays the contents of the original bag—showing that the IP addresses that comprise the set are identical to those in the bag.

Example 4.19: Creating an IP Set from a Bag with `rwbagtool --coverset`

```

<1>$ rwbagtool outgoingweb.bag --coverset \
  --output-path=outgoingweb.set
<2>$ rwsetcat outgoingweb.set \
| head -n 3
192.168.20.59
192.168.40.24
192.168.40.91
<3>$ rwbagcat outgoingweb.bag \
| head -n 3
 192.168.20.59|           1392516|
 192.168.40.24|           877782|
 192.168.40.91|           124548|

```

**Hint:** Be careful of bag contents when using `rwbagtool --coverset`. Since `rwbagtool` does not limit operations by the type of keys contained within a bag, the `--coverset` parameter will interpret the keys as IP addresses even if they are actually protocol or port keys. This will lead to confusion and analysis errors!

### Intersecting Bags and IPsets

You may want to find out whether the IP addresses in an IPset are also contained in a bag. You may also want to limit the contents of a bag to a specific group of IP addresses—for instance, those on a specific subnet.

The `rwbagtool --intersect` and `--complement-intersect` parameters are used to intersect an IPset with a bag. Example 4.20 shows how to use these parameters to extract a specific subnet.

Example 4.20: Using `rwbagtool --intersect` to Extract a Subnet

```
<1>$ echo '10.0.20.x' >f.set.txt
<2>$ rwsetbuild f.set.txt f.set
<3>$ rwbagtool x.bag --intersect=f.set --output-path=xf.bag
<4>$ rwbagcat x.bag
  10.0.20.58|          522|
  10.0.20.59|        1652|
  67.215.0.55|         88|
 117.34.28.84|         12|
   155.6.3.10|         30|
   155.6.4.10|         30|
192.168.200.10|       3913|
<5>$ rwbagcat xf.bag
  10.0.20.58|          522|
  10.0.20.59|        1652|
```

## 4.2.6 Masking IP Addresses

When working with IP addresses and utilities such as `rwuniq` and `rwstats`, you may want to analyze activity across *networks* rather than individual IP addresses. For example, you may wish to examine all of the activity originating from the /24s constituting the enterprise network rather than generating an individual entry for each address. To perform this type of analysis, use the `rwnetmask` command to reduce IP addresses to prefix values of a parameterized length. See Appendix C.9 for more information about `rwnetmask` or type `rwnetmask --help` at the command line.

The query in Example 4.21 is followed by an `rwnetmask` call to retain only the first 24 bits (three octets) of source IPv4 addresses, as shown by the `rwcut` output.

Example 4.21: Abstracting Source IPv4 addresses with `rwnetmask`

```
<1>$ rfilter --type=out,outweb --start-date=2015/06/02 \
  --end-date=2015/06/18 --sensors=S0,S1 --protocol=6 \
  --max-pass-records=3 --pass=stdout \
| rwnetmask --4sip-prefix-length=24 --4dip-prefix-length=24 \
| rwcut --fields=1-5
      sIP|          dIP|sPort|dPort|pro|
  10.0.40.0| 192.168.124.0| 1065| 591| 6|
  10.0.40.0| 192.168.166.0| 1066| 591| 6|
  10.0.40.0| 192.168.40.0|58083| 88| 6|
```

As Example 4.21 shows, `rwnetmask` replaces the last 8 bits<sup>8</sup> of the source and destination IP addresses with zero, so all IP addresses in the 10.0.40/24 block (for example) will be masked to produce the same IP address.

<sup>8</sup>32 bits total for an IPv4 address minus the 24 bits specified in the command for the prefix length leaves 8 bits to be masked.

This replaces both source and destination IP addresses with zero. With `rwnetmask`, an analyst can use any of the standard SiLK utilities on networks in the same way the analyst would use the utilities on individual IP addresses.

## 4.2.7 Working With IPsets

Iterative multi-path analyses commonly result in multiple SiLK IPset files. These files are usually named to describe their association with a specific aspect of analysis, such as byte thresholds as discussed in Section 2.2.8. As analyses progress, however, it is necessary to understand how IPsets compare and contrast. `rwsetbuild`, `rwsettool`, and `rwsetmember` are three important SiLK IPset tools that are often used together to identify network infrastructure and traffic flow.

### Creating IPsets from Text Files

`rwsetbuild` creates binary IPsets from text input files. It can be used to build reference sets to start an analysis, as previously discussed in Section 2.2.8. See Appendix C.17 for a command summary or type `rwsetbuild --help` at the command line.

Example 4.22 shows how to build an IPset of the FCCX-15 internal network reference from the `ipblocks` statements in the `sensors.conf` configuration file. Command 1 displays the first five lines of the `monitored_nets.set.txt` textual input file. Command 2 shows the use of the `rwsetbuild` to build the binary IPset from text. Finally, command 3 shows the file command to verify the `monitored_nets.set` filetype.

Example 4.22: Generating a Monitored Address Space IPset with `rwsetbuild`

---

```
<1>$ head -n 5 monitored_nets.set.txt
# Text file of monitored networks
# Build from sensor.conf ipblocks
10.0.10.0/24
10.0.20.0/24
10.0.30.0/24
<2>$ rwsetbuild monitored_nets.set.txt monitored_nets.set
<3>$ file monitored_nets.set
monitored_nets.set: SiLK, IPSET v2, Little Endian, LZ0 compression
```

---

Example 4.23 shows a similar approach to build an IPset of the broadcast address space with `rwsetbuild`.

Example 4.23: Generating a Broadcast Address Space IPset with `rwsetbuild`

---

```
<1>$ echo 255.255.255.255 >broadcast.set.txt
<2>$ rwsetbuild broadcast.set.txt broadcast.set
<3>$ file broadcast.set
broadcast.set: SiLK, IPSET v2, Little Endian, LZ0 compression
```

---

### Manipulating IPsets: DNS Server Example

Once you have constructed SiLK IPsets, use the `rwsettool` command for manipulating them. It provides common algebraic set operations for arbitrary numbers of IPset files. See Appendix C.16 for a summary of its syntax and most of its parameters or type `rwsettool --help` at the command line.

Example 4.24 shows an example of how to combine two sets to create a comprehensive IPset of the FCCX-15 internal network. It uses the `rwsettool --union` operation to combine both input set files, resulting in a summary set file, `internal_nets.set`. This file represents the FCCX-15 internal network addresses, including IP broadcasts that should not route to public IP space.

Example 4.24: Performing an IPset Union with `rwsettool`

---

```
<1>$ rwsettool --union monitored_nets.set broadcast.set \
>internal_nets.set
```

---

Analysts should determine the time period required for an analysis after creating reference sets. Example 4.25 shows how to use `rwsiteinfo` with the `--fields=repo-start-date,repo-end-date` parameter to determine the full time range of FCCX-15 data: 2015/06/02T13:00:00 to 2015/06/18T18:00:00.

Example 4.25: Displaying Repository Dates with `rwsiteinfo`

---

```
<1>$ rwsiteinfo --fields=repo-start-date,repo-end-date
      Start-Date |           End-Date |
2015/06/02T13:00:00 | 2015/06/18T18:00:00 |
```

---

The dates identified in Example 4.25 are then used as input to the `rwfilter` command to inventory all out type Domain Name System (DNS) servers, as shown in Example 4.26. Command 1 shows how to use `rwfilter` to query the entire FCCX-15 repository for DNS servers, which carry out type traffic on port 53 (`--dport=53`) with the UDP protocol (`--protocol=17`). It saves those IP addresses to the `dns_servers_out.set` file. Command 2 shows that there were outbound requests to 22 DNS servers on port 53/UDP during the period 2015/06/02 to 2015/06/18.

Example 4.26: Counting Outbound DNS Servers with `rwset`

---

```
<1>$ rwfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --protocol=17 --type=out --dport=53 --pass=stdout \
| rwset --dip-file=dns_servers_out.set
<2>$ rwsetcat --count dns_servers_out.set
22
```

---

Although Example 4.26 generated a comprehensive IPset of outbound DNS servers (`dns_servers_out.set`), it contains all servers that carry out type traffic. This means that DNS servers that are contained within the network perimeter may also reside in the resulting IPset. To differentiate between the internal and external network, IP addresses of internal DNS servers must be removed from the set.

Example 4.27 shows how to use the `rwsettool --difference` option to remove the IP addresses of internal DNS servers from the set of DNS servers that reside outside the network perimeter. Command 1 shows how to create the `external_dns_servers.set` set file by finding the difference between the DNS servers

contained in `dns_servers_out.set`, but not contained in `internal_nets.set`. Command 2 shows a total of 16 DNS servers that reside external to the network in the FCCX-15 data.

Example 4.27: Finding IPset Differences with `rwsettool`

---

```
<1>$ rwsettool --difference dns_servers_out.set \
    internal_nets.set >external_dns_servers.set
<2>$ rwsetcat --count external_dns_servers.set
16
```

---

The remaining internal DNS servers can be identified with the `rwsettool --symmetric-difference` option. A symmetric difference is the elements of two sets that are members of either set, but not members of both. Example 4.28 shows how the `internal_dns_server.set` set file is generated from a symmetric difference of the `external_dns_servers.set` and `dns_servers_out.set` files, finding a total of 6 internal DNS servers.

Example 4.28: Finding IPset Symmetric Difference with `rwsettool`

---

```
<1>$ rwsettool --symmetric-difference external_dns_servers.set \
    dns_servers_out.set >internal_dns_servers.set
<2>$ rwsetcat --count internal_dns_servers.set
6
```

---

## Using Set Membership to Understand Traffic Flow

Once DNS server infrastructure is identified, multi-path analyses may also require identifying traffic flow to specific DNS servers. Use the `rwsetmember` command to identify if an IP address or pattern is contained within one or more IPset files. This command shows how network traffic flows through a network infrastructure.

`rwsetmember` begins by building per-sensor IPset inventories of outbound traffic to DNS servers, as shown in Example 4.29.

Command 1 shows how to use `rwsiteinfo` to generate a list of sensors for the FCCX-15 dataset. This list of sensors is then used in commands 2-3 of the `dns_servers_by_sensor.sh` script to loop through each sensor name and build an IPset of the DNS servers that are monitored by each sensor.

After creating the per-sensor DNS server IPsets, we can use the `rwsetmember` command as shown in Example 4.30 to identify sensors that monitor specific external DNS servers.

- Command 1 shows how to use `rwsetmember` to identify the sensors that monitor 8.8.x.x DNS servers. The results indicate that sensors `S0`, `S1`, `S2`, `S3`, and `S12` logged DNS requests to 8.8.x.x IP addresses.
- Command 2 shows how these sensors can be combined into a list with `rwsiteinfo` to display their descriptions. The output from the `rwsiteinfo` command shows that DNS clients in the `Div0Ext`, `Div1Ext`, `Div0Int`, `Div1Int1`, and `Div1svc` monitored networks execute DNS queries to 8.8.x.x internet servers.

Example 4.29: Grouping Outbound DNS Servers by Sensor

```

<1>$ rwsiteinfo --fields=sensor:list
                                                    Sensor:list|
S0,S1,S2,S3,S4,S5,S6,S7,S8,S9,S10,S11,S12,S13,S14,S15,S16,S17,S18,S19,S20,S21|
<2>$ cat dns_servers_by_sensor.sh
SDATE="2015/06/02"
EDATE="2015/06/18"
SENSORS="S0 S1 S2 S3 S4 S5 S6 S7 S8 S9 S10 S11 S12 S13 S14"
SENSORS+=" S15 S16 S17 S18 S19 S20 S21"
for SENSOR in $SENSORS; do
    rwfilter --start-date=$SDATE --end-date=$EDATE --type=out \
    --proto=17 --dport=53 --sensor=$SENSOR --pass=stdout \
    | rwset --dip-file="$SENSOR"_dns_servers_out.set
done
<3>$ sh dns_servers_by_sensor.sh

```

Example 4.30: Identifying DNS Traffic Flow

```

<1>$ rwsetmember 8.8.x.x S*.set
S0_dns_servers_out.set
S12_dns_servers_out.set
S1_dns_servers_out.set
S2_dns_servers_out.set
S3_dns_servers_out.set
<2>$ rwsiteinfo --fields=sensor,describe-sensor \
--sensors=S0,S1,S2,S3,S12
Sensor|Sensor-Description|
S0|Div0Ext|
S1|Div1Ext|
S2|Div0Int|
S3|Div1Int1|
S12|Div1svc|

```

The `rwsettool --intersection` option can also be used to identify infrastructure that shares traffic flows. This option intersects IP addresses that are members of two sets. Command 1 of Example 4.31 shows using `rwsettool` with `--intersection` to identify the DNS servers that both the S0 and S1 sensors monitor.

Example 4.31: Identifying Shared DNS Monitoring

---

```
<1>$ rwsettool --intersect S0_dns_servers_out.set \
    S1_dns_servers_out.set | rwsetcat
8.8.4.4
8.8.8.8
67.215.0.5
128.8.10.90
128.63.2.53
192.5.5.241
192.33.4.12
192.36.148.17
192.58.128.30
192.112.36.4
192.203.230.10
192.228.79.201
193.0.14.129
198.41.0.4
199.7.83.42
202.12.27.33
```

---

### Displaying the Contents of IPsets

Use the `rwsetcat` command to view the contents of IPsets. It reads one or more set files, then displays the IP addresses in each file or prints out statistics about the set in each file. See Appendix C.15 for a command summary or type `rwsetcat --help` at the command line.

In Example 4.32, the call to `rwsetcat` prints out all the addresses in the set; IP addresses appear in ascending order.

Example 4.32: Displaying the Contents of IP Sets with `rwsetcat`

---

```
<1>$ rwsetcat medtcp-dest.set | head -n 5
192.168.45.27
192.168.61.26
```

---

In addition to printing out IP addresses, `rwsetcat` can also perform counting and statistical reporting, as shown in Example 4.33. These features are useful for describing the set without dumping out all the IP addresses in the set. Since sets can have any number of addresses, counting with `rwsetcat` tends to be much faster than counting via text tools such as `wc`.

- Command 1 shows how many IP addresses are in the IPset.
- Command 2 shows summary statistics and network structure for the addresses in the IPset.



Example 4.33 also shows the wide variety of network information that can be displayed by using the `--network-structure` parameter.

- In Command 3, there are no *list-lengths* and no *summary-lengths*. As a result, a default array of summary lengths is supplied.
- In Command 4 there is a *list-length*, but no slash (/) introducing *summary-lengths*, so the netblock with the specified prefix length is listed, but no summary is produced.
- In Command 5, a prefix length is supplied that is sufficiently large to list multiple netblocks.
- Command 6 shows two prefix lengths in *list-lengths*.
- Command 7 shows that a prefix length of zero (no network bits, so no choice of networks) treats the entire address space as a single network and labels it **TOTAL**.
- Command 8 shows that summarization occurs not only for the *summary-lengths* but also for every prefix length in *list-lengths* that is larger than the current list length.
- In Command 9, the slash introduces *summary-lengths*, but the array of summary lengths is empty; as a result, the word “hosts” appears as if there will be summaries, but there aren’t any.
- In Command 10, the **S** replaces the slash and summary lengths, so default summary lengths are used.
- In Command 11, the list length is larger than the smallest default summary length, so that summary length does not appear.
- In Command 12, **H** (host) is used for a list length.
- Command 13 shows that **H** is equivalent to 32 for IPv4.

Example 4.33: `rwsetcat` Options for Showing Structure

```
<1>$ rwsetcat medtcp-dest.set --count-ips
93
<2>$ rwsetcat medtcp-dest.set --print-statistics
Network Summary
  minimumIP =      10.0.40.20
  maximumIP = 192.168.166.233
      93 hosts (/32s),    0.000002% of 2^32
      2 occupied /8s,    0.781250% of 2^8
      2 occupied /16s,   0.003052% of 2^16
      20 occupied /24s,  0.000119% of 2^24
      66 occupied /27s,  0.000049% of 2^27
<3>$ rwsetcat medtcp-dest.set --network-structure
TOTAL| 93 hosts in 2 /8s, 2 /16s, 20 /24s, and 66 /27s
<4>$ rwsetcat medtcp-dest.set --network-structure=4
  0.0.0.0/4| 10
 192.0.0.0/4| 83
<5>$ rwsetcat medtcp-dest.set --network-structure=18
 10.0.0.0/18| 10
 192.168.0.0/18| 15
 192.168.64.0/18| 27
 192.168.128.0/18| 41
```

```

<6>$ rwssetcat medtcp-dest.set --network-structure=4,18
 10.0.0.0/18 | 10
0.0.0.0/4 | 10
 192.168.0.0/18 | 15
 192.168.64.0/18 | 27
 192.168.128.0/18 | 41
192.0.0.0/4 | 83
<7>$ rwssetcat medtcp-dest.set --network-structure=0,18
 10.0.0.0/18 | 10
 192.168.0.0/18 | 15
 192.168.64.0/18 | 27
 192.168.128.0/18 | 41
TOTAL | 93
<8>$ rwssetcat medtcp-dest.set --network-structure=4,18/24
 10.0.0.0/18 | 10 hosts in 2 /24s
0.0.0.0/4 | 10 hosts in 1 /18 and 2 /24s
 192.168.0.0/18 | 15 hosts in 3 /24s
 192.168.64.0/18 | 27 hosts in 6 /24s
 192.168.128.0/18 | 41 hosts in 9 /24s
192.0.0.0/4 | 83 hosts in 3 /18s and 18 /24s
<9>$ rwssetcat medtcp-dest.set --network-structure=4/
 0.0.0.0/4| 10 hosts
 192.0.0.0/4| 83 hosts
<10>$ rwssetcat medtcp-dest.set --network-structure=4S
 0.0.0.0/4| 10 hosts in 1 /8, 1 /16, 2 /24s, and 4 /27s
 192.0.0.0/4| 83 hosts in 1 /8, 1 /16, 18 /24s, and 62 /27s
<11>$ rwssetcat medtcp-dest.set --network-structure=12S
 10.0.0.0/12| 10 hosts in 1 /16, 2 /24s, and 4 /27s
 192.160.0.0/12| 83 hosts in 1 /16, 18 /24s, and 62 /27s
<12>$ rwssetcat medtcp-dest.set --network-structure=H \
| head -n 5
 10.0.40.20|
 10.0.40.23|
 10.0.40.53|
 10.0.40.83|
 10.0.50.11|
<13>$ rwssetcat medtcp-dest.set --network-structure=32 \
| head -n 5
 10.0.40.20|
 10.0.40.23|
 10.0.40.53|
 10.0.40.83|
 10.0.50.11|

```

## 4.2.8 Indicating Flow Relationships

A useful step in a multi-path analysis is to identify a set of flow records that have common attributes—for instance, records that are part of the same TCP session. Use `rwgroup` and `rwmatch` to label a set of flow records that share attributes. This identifier, or group ID, is stored in the next-hop IP (`nhIP`) field. It can be manipulated as an IP address (that is, either by directly specifying a group ID or by using IPsets). The two tools generate group IDs in different ways.

- `rwgroup` scans a file of flow records and groups records with common attributes, such as source or destination IP address pairs.
- `rwmatch` groups records of different types (typically, incoming and outgoing types), creating a file containing groups that represent TCP sessions or groups that represent other behavior.

**Hint:** To improve scalability, the grouping tools require the data they process to first be sorted using `rwsort`. The sorted data must be sorted on the criteria fields: in the case of `rwgroup`, the ID field and delta fields; in the case of `rwmatch`, start time and the fields specified in the `--relate` parameter(s).

### Labeling Flow Records Based on Common Attributes

The `rwgroup` command groups flow records that have common field values. Grouped records can be output separately (with each record in the group having a common ID) or summarized by a single record. Applications of `rwgroup` include the following:

- grouping together all flow records for a long-lived session: By specifying that records are grouped together by their port numbers and IP addresses, an analyst can assign a common ID to all the flow records making up a long-lived session.
- reconstructing web sessions: Due to diversified hosting and caching services such as Akamai<sup>®</sup>, a single webpage on a commercial website is usually hosted on multiple servers. For example, the images may be on one server, the HTML text on a second server, advertising images on a third server, and multimedia on a fourth server. An analyst can use `rwgroup` to tag web traffic flow records from a single user that are closely related in time and then use that information to identify individual webpage fetches.
- counting conversations: An analyst can group all the communications between two IP addresses together and see how much data was transferred between both sites regardless of port numbers. This is particularly useful when one site is using a large number of ephemeral ports.

See Appendix C.26 for a command summary or type `rwgroup --help` at the command line.

Flow records are grouped when the fields specified by `--id-fields` are identical and the field specified by `--delta-field` matches within a value less than or equal to the value specified by `--delta-value`.

Creating a group is a two-step process:

1. Sort the records using `rwsort` as described in Section 2.2.7. `rwgroup` requires input records to be sorted by the fields specified in `--id-fields` and `--delta-field`.
2. Run `rwgroup` to create a group that matches the grouping criteria specified by `--id-fields`, `--delta-field`, and `--delta-value`. Records in the same group are assigned a common group ID.

`rwgroup` outputs a stream of flow records. Each record's next-hop IP address field is set to the value of the group ID.

**Grouping Records By Session.** The most basic use of `rwgroup` is to group together flow records that constitute a single longer session, such as the components of a single FTP session (or, in the case of Example 4.34, a Microsoft<sup>®</sup> Distributed File System Replication Service session). To do this, the example does the following:

1. Command 1 uses `rwfilter` to pull the desired flow records from the repository. It then uses the `rwsort` command to sort these records by source and destination IP address, source and destination port, and start time.
2. Command 2 uses `rwgroup` to group together flow records that have closely-related start times.
3. Command 3 uses the `rwfilter` and `rwcut` commands to display grouped records. It filters for records with `--next-hop-id=0.0.0.1,28` (the group IDs), then displays the source and destination IP addresses, source and destination ports, and the group ID (nhIP).

This creates a group of records that comprise a session.

Example 4.34: Grouping Flows of a Long Session with `rwgroup`

```
<1>$ rwfilter --type=in,out --start-date=2015/06/02 \
--end-date=2015/06/18 --packets=4- --protocol=6 \
--bytes-per-packet=60- --duration=1000- --pass=stdout \
| rwsort --fields=1,2,3,4,sTime --output-path=sorted.rw
<2>$ rwgroup sorted.rw --id-fields=1,2,3,4 --delta-field=sTime \
--delta-value=3600 --output-path=grouped.rw
<3>$ rwfilter grouped.rw --next-hop-id=0.0.0.1,28 --pass=stdout \
| rwcut --fields=1-4,nhIP
```

sIP	dIP	sPort	dPort	nhIP
10.0.40.20	192.168.40.20	55425	5722	0.0.0.1
10.0.40.20	192.168.40.20	55425	5722	0.0.0.1
10.0.40.20	192.168.40.20	55425	5722	0.0.0.1
192.168.40.20	10.0.40.20	5722	55425	0.0.0.28
192.168.40.20	10.0.40.20	5722	55425	0.0.0.28
192.168.40.20	10.0.40.20	5722	55425	0.0.0.28

**Thresholding Groups By Number of Records.** By default, `rwgroup` outputs one flow record for every flow record it receives as input. You can set a threshold for flow record output by using the `--rec-threshold` parameter, as shown in Example 4.35. This parameter specifies that `rwgroup` only passes records that belong to a group with at least as many records as given in `--rec-threshold`. All other records are dropped silently.

This allows you to filter out smaller groups of records. For instance, if you are only interested in grouping significant amounts of traffic, you could drop groups with low flow counts. Example 4.35 shows how this thresholding works. In the first case, there are several low-flow-count groups. When `rwgroup` is invoked with `--rec-threshold=4`, these groups are discarded by `rwgroup`, while the groups with 4 or more flow records are output.

**Generating a Single Summary Record for a Group.** `rwgroup` can also generate a single summary record with the `--summarize` parameter. When this parameter is used, `rwgroup` only produces a single record for each group. The summary record uses the first record in the group for its addressing information (IP addresses, ports, and protocol). The total number of bytes and packets for the group is recorded in the summary record's corresponding fields, and the start and end times for the record will be the extrema for that group.<sup>9</sup>

<sup>9</sup>This is only a quick version of condensing long flows—TCP flags, termination conditions, and application labeling may not be properly reflected in the output.

Example 4.35: Dropping Trivial Groups with `rwgroup --rec-threshold`

```

<1>$ rwgroup sorted.rw --id-fields=1,2,3,4 --delta-field=sTime \
  --delta-value=3600 \
| rwcut --num-recs=10 --field=1-5,nhIP
      sIP|          dIP|sPort|dPort|pro|          nhIP|
10.0.40.20| 192.168.40.20| 5722|60309| 6|          0.0.0.0|
10.0.40.20| 192.168.40.20|55425| 5722| 6|          0.0.0.1|
10.0.40.20| 192.168.40.20|55425| 5722| 6|          0.0.0.1|
10.0.40.20| 192.168.40.20|55425| 5722| 6|          0.0.0.1|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
<2>$ rwgroup sorted.rw --id-fields=1,2,3,4 --delta-field=sTime \
  --delta-value=3600 --rec-threshold=4 \
| rwcut --num-recs=10 --field=1-5,nhIP
      sIP|          dIP|sPort|dPort|pro|          nhIP|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|
10.0.50.12| 192.168.40.100| 3088| 8005| 6|          0.0.0.2|

```

Example 4.36 shows how summarizing works: The 10 original records are reduced to two group summaries, and the byte totals for those records are equal to the sum of the byte values of all the records in the group.

Example 4.36: Summarizing Groups with `rwgroup --summarize`

```
<1>$ rwgroup sorted.rw --id-fields=1,2,3,4 --delta-field=sTime \
--delta-value=3600 --rec-threshold=3 \
| rwcut --fields=1-5,bytes,nhIP --num-recs=10
  sIP|          dIP|sPort|dPort|pro|      bytes|      nhIP|
 10.0.40.20| 192.168.40.20|55425| 5722| 6|      5523| 0.0.0.1|
 10.0.40.20| 192.168.40.20|55425| 5722| 6|     21084| 0.0.0.1|
 10.0.40.20| 192.168.40.20|55425| 5722| 6|     11500| 0.0.0.1|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    977689| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|    959308| 0.0.0.2|
<2>$ rwgroup sorted.rw --id-fields=1,2,3,4 --delta-field=sTime \
--delta-value=3600 --rec-threshold=3 --summarize \
| rwcut --fields=1-5,bytes,nhIP --num-recs=5
  sIP|          dIP|sPort|dPort|pro|      bytes|      nhIP|
 10.0.40.20| 192.168.40.20|55425| 5722| 6|     38107| 0.0.0.1|
 10.0.50.12| 192.168.40.100| 3088| 8005| 6|   46529266| 0.0.0.2|
 10.0.50.12| 192.168.40.100| 3089| 8005| 6|   46726003| 0.0.0.3|
 10.0.50.12| 192.168.40.100| 3090| 8005| 6|   46497279| 0.0.0.4|
 10.0.50.12| 192.168.40.100| 3091| 8005| 6|   46448588| 0.0.0.5|
```

**Grouping Records via IPsets.** For any data file, calling `rwgroup` with the same `--id-fields` and `--delta-field` values will result in the same group IDs being assigned to the same records. As a result, an analyst can use `rwgroup` to manipulate groups of flow records where the group has a specific attribute. This can be done by using `rwgroup` and IPsets, as shown in Example 4.37.

1. Command 1 uses `rwfilter` to filter for traffic with the TCP protocol (`--protocol=6`) and destination ports 20 and 21 (`--dport=20,21`). It then calls `rwsort` to sort the data and uses `rwgroup` to convert the results into a file, `out.rw`, grouped as FTP communications between two sites. All TCP port 20 and 21 communications between two sites are part of the same group.
2. Command 2 filters through the collection of groups for those group IDs (as next-hop IP addresses stored in `control.set`) that use FTP control.
3. Finally, Command 3 uses that next-hop IPset to pull out all the flows (ports 20 and 21) in groups that had FTP control (port 21) flows.

## Labeling Matched Groups of Flow Records

As part of a larger analysis, you may want to group network flow records. For instance, you could group records from both sides of a bidirectional session, such as HTTP requests and responses, for further exam-

Example 4.37: Using `rwgroup` to Identify Specific Sessions

```

<1>$ rfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --protocol=6 --type=out --dport=20,21 --pass=stdout \
| rwsort --fields=1,2,sTime \
| rwgroup --id-fields=1,2 --output-path=out.rw
<2>$ rfilter out.rw --dport=21 --pass=stdout \
| rwset --nhip-file=control.set
<3>$ rfilter out.rw --nhipset=control.set --pass=stdout \
| rwcut --fields=1-5,sTime --num-recs=5

```

sIP	dIP	sPort	dPort	pro	sTime
192.168.70.10	10.0.40.83	65360	21	6	2015/06/16T20:38:57.018
192.168.70.10	10.0.40.83	65360	21	6	2015/06/16T20:38:57.018
192.168.70.10	10.0.40.83	65360	21	6	2015/06/16T20:38:57.146
192.168.70.10	10.0.40.83	57096	21	6	2015/06/17T04:41:35.525
192.168.70.10	10.0.40.83	57096	21	6	2015/06/17T04:41:35.525

ination. You may also wish to create groups with more flexible matching, such as matching groups across protocols to identify `traceroute` messages, which use the UDP and ICMP protocols.

Use the `rwmatch` to create *matched groups*. A matched group consists of an initial record (usually a *query*) followed by one or more *responses* and (optionally) additional queries. (For more information about this command's syntax and common options, see Appendix C.25 or type `rwmatch --help` at the command line.)

A response is a record that is related to the query (as specified in the `rwmatch` command). However, it is collected from a different direction or from a different router. As a result, the fields relating the two records may be different. For example, the source IP address in one record may match the destination IP address in another record.

A relationship in `rwmatch` is established using the `--relate` parameter, which takes two numeric field IDs separated by a comma (e.g., `--relate=3,4` or `--relate=5,5`). The first value corresponds to the field ID in the query file,. The second value corresponds to the field ID in the response file. For example, `--relate=1,2` states that the source IP address in the query file must match the destination IP address in the response file. The `rwmatch` tool will process multiple relationships, but each field in the query file can be related to, at most, one field in the response file.

The two input files to `rwmatch` must be sorted before matching. The same information provided in the `--relate` parameters, plus `sTime`, must be used for sorting. The first fields in the `--relate` value pairs, plus `sTime`, constitute the sort fields for the query file. The second fields in the `--relate` value pairs, plus `sTime`, constitute the sort fields for the response file.

The `--relate` parameter always specifies a relationship from the query to the responses, so specifying `--relate=1,2` means that the records match if the source IP address in the query record matches the destination IP address in the response. Consequently, when working with a protocol where there are implicit relationships between the queries and responses, especially TCP, these relationships must be fully specified. Example 4.38 shows the impact that not specifying all the fields has on TCP data. Note that the match relationship specified (query's source IP address matches response's destination IP address) results in all the records in the response matching the initial query record, even though the source ports in the query file may differ from a response's destination port (as seen with the third matched record).

Example 4.39 shows the relationships that could be specified when working with TCP or UDP. This example

Example 4.38: Using `rwmatch` with Incomplete Relate Values

```

<1>$ rfilter --type=in,out --start-date=2015/06/02 \
  --end-date=2015/06/18 --protocol=6 --dport=88 \
  --pass=stdout \
| rwsort --fields=1 --output-path=query.rw
<2>$ rfilter --type=in,out --start-date=2015/06/02 \
  --end-date=2015/06/18 --protocol=6 --sport=88 \
  --pass=stdout \
| rwsort --fields=2 --output-path=response.rw
<3>$ rwcut query.rw --fields=1-4,sTime --num-recs=4
      sIP|                dIP|sPort|dPort|                sTime|
10.0.40.26| 192.168.40.20|57288| 88|2015/06/17T17:06:38.871|
10.0.40.26| 192.168.40.20|57287| 88|2015/06/17T17:06:38.865|
10.0.40.26| 192.168.40.20|52176| 88|2015/06/16T16:22:08.659|
10.0.40.26| 192.168.40.20|57287| 88|2015/06/17T17:06:38.856|
<4>$ rwcut response.rw --fields=1-4,sTime --num-recs=4
      sIP|                dIP|sPort|dPort|                sTime|
192.168.40.20| 10.0.40.26| 88|57390|2015/06/17T17:32:21.501|
192.168.40.20| 10.0.40.26| 88|52724|2015/06/16T18:54:21.810|
192.168.40.20| 10.0.40.26| 88|52725|2015/06/16T18:54:21.818|
192.168.40.20| 10.0.40.26| 88|57389|2015/06/17T17:32:21.493|
<5>$ rwmatch --relate=1,2 query.rw response.rw stdout \
| rwcut --fields=1-4,nhIP --num-recs=10
      sIP|                dIP|sPort|dPort|                nhIP|
10.0.40.26| 192.168.40.20|57390| 88|                0.0.0.1|
192.168.40.20| 10.0.40.26| 88|57390|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|52724|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|52725|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|57389|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|51466|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|57321|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|57320|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|52030|                255.0.0.1|
192.168.40.20| 10.0.40.26| 88|52031|                255.0.0.1|

```



specifies a relationship between the query's source IP address and the response's destination IP address, the query's source port and the response's destination port, and the reflexive relationships between query and response.

`rwmatch` is designed to handle not just protocols where source and destination are closely associated, but also where partial associations are significant.

Example 4.39: Using `rwmatch` with Full TCP Fields

```
<1>$ rwsort query.rw --fields=1,2,sTime \
--output-path=squery.rw
<2>$ rwsort response.rw --fields=2,1,sTime \
--output-path=sresponse.rw
<3>$ rwmatch --relate=1,2 --relate=2,1 --relate=3,4 --relate=4,3 \
squery.rw sresponse.rw stdout \
| rwcut --fields=1-4,nhIP --num-recs=5
      sIP |          dIP | sPort | dPort |          nhIP |
10.0.40.26 | 192.168.40.20 | 52396 | 88 | 0.0.0.1 |
192.168.40.20 | 10.0.40.26 | 88 | 52396 | 255.0.0.1 |
10.0.40.26 | 192.168.40.20 | 51466 | 88 | 0.0.0.2 |
192.168.40.20 | 10.0.40.26 | 88 | 51466 | 255.0.0.2 |
10.0.40.26 | 192.168.40.20 | 51466 | 88 | 0.0.0.3 |
```

To establish a match group, a response record must first match a query record. For that to happen all the fields of the query record specified as first values in relate pairs must match all the fields of the response record specified as second values in the relate pairs. Additionally, the start time of the response record must fall in the interval between the query record's start time and its end time extended by the value of `--time-delta`. Alternatively, if the `--symmetric-delta` parameter is specified, the query record's start time may fall in the interval between the response record's start time and its end time extended by `--time-delta`. The record whose start time is earlier becomes the base record for further matching.

Additional target records from either file may be added to the match group. If the base and target records come from different files, field matching is performed with the two fields specified for each relate pair. If the base and target records come from the same file, field matching is done with the same field for both records.

In addition to matching the relate pairs, the target record's start time must fall within a time window beginning at the start time of the base record. If `--absolute-delta` is specified, the window ends at the base record's end time extended by `--time-delta`. If `--relative-delta` is specified, the window ends `--time-delta` seconds after the maximum end time of any record in the match group so far. If `--infinite-delta` is specified, time matching is not performed.

As with `rwgroup`, `rwmatch` sets the next-hop IP address field to an identifier common to all related flow records. However, `rwmatch` groups records from two distinct files into single groups. To indicate the origin of a record, `rwmatch` uses different values in the next-hop IP address field. Query records will have an IPv4 address where the first octet is set to zero; in response records, the first octet is set to 255. `rwmatch` *only* outputs queries that have a response grouped with all the responses to that query.

`rwmatch` discards queries that do not have a response and responses that do not have a query unless `--unmatched` is specified. It tells `rwmatch` to write unmatched query and/or response records to an output file or standard output.

- `q` writes query records. Unmatched query records have their next hop IPset to 0.0.0.0.

- **r** writes response records. Unmatched response records have their next hop IPset to 255.0.0.0.
- **b** writes both query and response records.

As a result, `rwmatch`'s output usually has fewer records than the total of the two source files. `rwgroup` can be used to compensate for this by merging all the query records for a single session into one record.

Example 4.39 matches all addresses and ports in both directions. As with `rwgroup`, `rwmatch` requires sorted data, and in the case of `rwmatch`, there is always an implicit time-based relationship controlled using the `--time-delta` parameter. As a consequence, *always* sort `rwmatch` data on the start time. (Example 4.38 generated the query and response files from a query that might not produce sorted records; Example 4.39 corrected this by sorting the input files before calling `rwmatch`.)

### 4.2.9 Managing IPset, Bag, and Prefix Map Files

During a multi-path analysis, an analyst typically performs many intermediate steps while isolating the behavior of interest. The `rwfileinfo` command prints information about binary SiLK files, helping you to manage the files generated as part of this process.

`rwfileinfo` can display information about all types of binary SiLK files: flow record files, IPset files, bag files, and prefix map (or pmap) files. Section 2.2.4 describes how to use `rwfileinfo` for viewing information about flow record files. The current section describes how to use `rwfileinfo` to view information about set, bag, and pmap files, such as the record count, file size, and the SiLK command(s) that created the file. Additionally, `rwfileinfo` displays information specific to each type of file. This information is displayed by default along with the rest of the file information, or can be explicitly viewed by using the `--fields` parameter.

- For IPset files, it shows the nodes, branches and leaves of the IPset (`--fields=ipset`).
- For bag files, it shows the key-value pairs that make up the bag (`--fields=bag`).
- For pmap files, it shows the internal prefix map name (`--fields=prefix-map`). If a prefix map was created without a map name, `rwfileinfo` returns an empty result for the prefix-map-specific field.

Example 4.40 show how `rwfileinfo` handles these files.

1. Commands 1, 2 and 3 create a set, a bag, and a pmap for the example.
2. Command 4 shows a full `rwfileinfo` result for the set file.
3. Commands 5, 6 and 7 show just the specific information for the IPset, bag, and pmap file, respectively.

---

 Example 4.40: `rwfileinfo` for Sets, Bags, and Prefix Maps
 

---

```

<1>$ rfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --sensors=S0,S1 --type=out,outweb --protocol=0- \
  --pass=stdout \
| rwbag --bag-file=sIPv4,flows,internal.bag
<2>$ rwbagtool --coverset --ipset-record-version=4 \
  --output-path=internal.set internal.bag
<3>$ rwpmapbuild --input-file=internal.pmap.txt \
  --output-file=internal.pmap
<4>$ rwfileinfo internal.set
internal.set:
  format(id)          FT_IPSET(0x1d)
  version              16
  byte-order           littleEndian
  compression(id)     lzo1x(2)
  header-length        56
  record-length        1
  record-version       4
  silk-version         3.16.1
  count-records        1184
  file-size            506
  ipset                IPv4
<5>$ rwfileinfo internal.set --fields=ipset
internal.set:
  ipset                IPv4
<6>$ rwfileinfo internal.bag --fields=bag
internal.bag:
  bag                  key: sIPv4 @ 4 octets; counter: records @ 8 octets
<7>$ rwfileinfo internal.pmap --fields=prefix-map
internal.pmap:
  prefix-map           v1: internal
  
```

---



## Chapter 5

# Case Studies: Intermediate Multi-path Analysis

This chapter features detailed case studies of multi-path analyses, using concepts from previous chapters. They employ the SiLK workflow, `rwfilter` and other SiLK tools, UNIX commands, and networking concepts to provide practical examples of multi-path analyses with network flow data.

Upon completion of this chapter you will be able to

- describe how to combine several tasks to form a multi-path analyses
- execute multi-path analyses with various SiLK tools in one automated program
- associate sets of IP addresses (IPsets) with network flow sensors

### 5.1 Building Inventories of Network Flow Sensors With IPsets

Flow sensors commonly monitor strategic points in enterprise networks where different network environments meet. This environmental complexity affects sensor flow collection and analyst knowledge as network infrastructure evolves. For example, multiple sensors may overlap their flow collection for failover purposes; as the network routes traffic, analysts may need to determine which sensor is the primary flow collector.

To mitigate these issues, analysts can create and maintain inventories of network sensors, making it easier to review and validate them. These sensor inventories consist of SiLK IPsets that contain internal network addresses monitored by a flow sensor. They are generated by applying the following multi-path analysis workflow.

1. Path 1 associates network addresses with a single sensor.
2. Path 2 associates network addresses of the remaining sensors.
3. Path 3 associates network shared addresses.
4. Finally, the results of each part of the multi-path analysis are merged to create a complete inventory of sensors.

This case is an example of how a multi-path analysis can be built from successive single-path analyses. It demonstrates how to build IPset inventories by merging the results of three single-path analyses into a single IPset inventory. Like the case study and command examples earlier in this guide, it uses the FCCX dataset described in Section 1.7. All analysis steps are shown in Example 5.1. Each part of the multi-path analysis is described in its own section.

Example 5.1: Building an IPset Inventory for Sensor S0

---

```

<1>$ rfilter --sensor=S0 --type=out,outweb --start=2015/06/01 \
  --end=2015/06/30 --proto=0- --pass=stdout \
| rset --sip-file=S0-out.set --copy=stdout \
| rwbag --bag-file=sipv4,flows,S0-out.bag
<2>$ rfilter \
  --sensor=S1,S2,S3,S4,S5,S6,S7,S8,S9,S10,S11,S12,S13,S14,S15,\
  S16,S17,S18,S19,S20,S21 \
  --type=out,outweb --start=2015/06/01 --end=2015/06/30 \
  --proto=0- --pass=stdout \
| rset --sip-file=S0-other.set --copy=stdout \
| rwbag --bag-file=sipv4,flows,S0-oth.bag
<3>$ rsettool --difference S0-out.set S0-other.set \
  --output=S0-only.set
<4>$ rwbagtool --compare=ge S0-out.bag S0-oth.bag \
  --output=S0-most.bag
<5>$ rwbagtool --coverset S0-most.bag --output=S0-most.set
<6>$ rsetcat --net=27 S0-most.set S0-only.set \
| cut -f1 -d '|' \
| rsetbuild stdin cand.set
<7>$ rfilter --sensor=S0 --type=in,inweb --start=2015/06/01 \
  --end=2015/06/30 --dipset=cand.set --pass=stdout \
| rset --dip-file=S0-in.set
<8>$ rsettool --difference S0-in.set S0-most.set S0-only.set \
  --output=S0-close.set
<9>$ rsettool --union S0-only.set S0-most.set S0-close.set \
  --output=S0.set

```

---

### 5.1.1 Path 1: Associate Addresses with a Single Sensor

To begin generating IPsets, you must select a sensor for the inventory (referenced as an *inventory sensor* for this case study).

1. Command 1 of Example 5.1 shows how to use the `rfilter` command to select all outbound traffic on inventory sensor `S0` from the repository.
2. Command 1 then calls the `rset` command to generate a source IPset of internal addresses (`S0-out.set`) that pass the `rfilter` query. This IPset contains all outbound source IP addresses monitored by sensor `S0` for the defined time period.
3. The `rset --copy` switch copies the `rfilter` results to `rwbag`, generating a SiLK bag of flow counts per source IP address for the same time period (`S0-out.bag`).

The resulting IPset and bag provide a detailed inventory of all source IP addresses and their flow volumes monitored by sensor `S0`.

### 5.1.2 Path 2: Associate Addresses of Remaining Sensors

Command 2 of Example 5.1 shows the next path in the multi-path analysis: associating the addresses for non-inventory sensors in the repository. This procedure is similar to that of command 1. However, command 2 selects outbound traffic for the remaining repository sensors (`S1` through `S21`) and generates source IP address sets and bag files for the same time period (`S0-other.set`, `S0-oth.bag`).

### 5.1.3 Path 3: Associate Shared Addresses

The third path of the multi-path analysis associates the addresses that may be monitored (or *shared*) by multiple sensors.

#### Identify Uniquely-Monitored Addresses

To begin path three, you need to identify addresses that are uniquely monitored by the inventory sensor, `S0`.

1. Command 3 of Example 5.1 shows how to use `rwsettool --difference` to generate an IPset of addresses that are uniquely monitored by sensor `S0` (`S0-only.set`). This command compares the set of addresses from command 1 that are monitored by `S0` (`S0-out.set`) to the set of addresses from command 2 that are monitored by all the other sensors (`S0-other.set`), and saves the IP addresses that are only found in `S0-out.set` to `S0-only.set`.
2. Command 4 of Example 5.1 checks to see if the inventory sensor `S0` is the primary flow collection sensor. It compares the IP address flow volumes of the inventory sensor to those of the other sensors to determine if `S0` monitors the majority of network traffic for an IP address. To perform this comparison, it uses `rwbagtool` to compare the bag files created in commands 1 and 2 (`S0-out.bag` and `S0-oth.bag`). IP addresses that are monitored by the inventory sensor and have flow volumes greater than or equal to those monitored by the non-inventory sensors are saved to the `S0-most.bag` file.
3. Command 5 shows how to use `rwbagtool --coverset` to generate an IPset (`S0-most.set`) from the `S0-most.bag` file.

#### Find Asymmetric and Missing Flow Data

To continue Path 3 of the multi-path analysis, you must account for asymmetric and missing network flow data. Causes of this include outages, routing, and network backdoors. To address each case, the third path should identify inbound network traffic to internal hosts closely adjacent to other inventory sensor addresses. This can be accomplished using a small CIDR range, such as `/27`.

1. Command 6 of Example 5.1 shows how to use the `rwsetcat` and `rwsetbuild` commands to build a `cand.set` candidate set from the `S0-only` and `S0-most` IPset files. The `--net=27` parameter specifies the CIDR range for closely adjacent addresses in the two sets.

2. This candidate set is then used as a destination IPset (`--dipset=cand.set`) for the `rwfilter` query in Command 7 to identify any remaining addresses not previously found in Commands 1 and 2.
3. Command 8 then uses `rwsettool --difference` to create the `S0-close.set` IPset, which contains IP addresses that are closely adjacent to inventory sensor addresses.

### 5.1.4 Merge Address Results

After completing the third path of the multi-path analysis, you will have three subsets of the sensor inventory: `S0-only.set`, `S0-most.set`, and `S0-close.set`. To complete the full sensor inventory, all three subsets must be merged together into a single IPset that contains the inventory of IP addresses associated with the inventory sensor, `S0`.

Command 9 of Example 5.1 shows how to generate the final sensor inventory with the `rwsettool --union` switch. This command performs an algebraic *union* on `S0-only.set`, `S0-most.set`, and `S0-close.set` to produce `S0.set`, which contains all of the IP addresses in the sensor inventory.

Analysts can use and maintain the sensor inventory IPset to gain a better understanding of the networks monitored by a specific network flow sensor (in this example, `S0`) and identify any situational awareness changes.

## 5.2 Automating IPset Inventories of Network Flow Sensors

Section 5.1 described how to manually generate a sensor IPset inventory by performing a multi-path analysis. This process is fine if you only need to inventory a single sensor. However, what if you need to inventory all of the sensors in the SiLK repository? Manually executing this workflow for every sensor would quickly become time consuming and represents an inefficient use of analyst time and resources.

Instead, create a shell program to automatically execute this workflow. An automated program can use the `rwsiteinfo` command to compile a list of *all* sensors in the repository, then execute the analytic from Section 5.1 for every sensor. Automating the workflow is more efficient and accurate than manually executing it for each sensor. It makes compiling and updating a complete sensor inventory much easier: simply run the program as needed to automatically generate a new inventory for all sensors that are currently in the SiLK repository. This improves your situational awareness of the network.

This section discusses Example 5.2, a simple Bash shell program that automates the process of building sensor IPset inventories. This program consists of two main sections, a *header* and *loop*, that are discussed in detail below.

### 5.2.1 Program Header

The first section of this automated program is the *header*, which contains information used throughout program execution.

Line 1 of Example 5.2, commonly known as the *she-bang*, identifies this as a *Bourne-again (Bash)* shell program to the system shell. The comment in Line 2 helps users to understand the overall purpose of the program.

Lines 5-8 define variables that are referenced throughout the remainder of the program.



- The `START` and `END` variables specify the dates that are used with `rwfilter --start-date` and `--end-date` switches.
- The `INBOUND` and `OUTBOUND` variables specify the options for inbound and outbound `rwfilter` queries that are executed in the program's *for loop*.

### 5.2.2 Program Loop

The remainder of the program, a Bash *for loop*, executes most of the program's automated tasks.

Line 11 of Example 5.2 loops through the sensors in the SiLK repository. The `rwsiteinfo --sensor` command generates a complete list of sensors in the repository. The *for loop* cycles through this list to generate an inventory for each sensor.

Lines 13-15 define a variable (`FILES` in this example) of temporary files that will be deleted at the end of the *for loop* (Line 47).

Line 17 follows with a visual output of the current inventory sensor using the UNIX `echo` command. This indicates that the program has started building an IPset inventory for that sensor.

With the exception of Lines 23-25, the remainder of the program repeats the analytic in Example 5.1 using the variables defined in the program header. Lines 23-25 are an addition to the original analytic. They call `rwsiteinfo` with the `--fields=sensor:list` option to generate a comma-delimited list of sensors. This list of sensors is then passed to a series of `sed` commands to remove the inventory sensor, effectively building an `OTHERS` variable that contains a list of non-inventory sensors. The `OTHERS` variable is then used to build the set of non-inventory sensors used in the rest of the analytic.

Line 49 again displays the name of the current sensor and indicates that the program has finished building an inventory for that sensor.

## Example 5.2: Automating IPset Inventories

```

1  #!/bin/bash
   # Script to automate sensor IPset inventories

4  # Variables
   START="2015/06/01"
   END="2015/06/30"
7  OUTBOUND="--type=out,outweb --start=$START --end=$END --proto=0-"
   INBOUND="--type=in,inweb --start=$START --end=$END --dipset=cand.set"

10 # Loop through each sensor in the repository
   for SEN in $(rwsiteinfo --no-titles --no-final --fields=sensor); do

13     FILES="$SEN-out.set $SEN-out.bag $SEN-close.set $SEN-other.set"
       FILES+=" $SEN-oth.bag $SEN-only.set cand.set"
       FILES+=" $SEN-most.set $SEN-in.set $SEN-most.bag"

16     echo "Sensor: $SEN -- Start"

19     rfilter --sensor=$SEN $OUTBOUND --pass=stdout \
       | rwsset --sip-file=${SEN}-out.set --copy=stdout \
       | rwbag --bag-file=sipv4,flows,${SEN}-out.bag

22     OTHERS=$(rwsiteinfo --no-titles --no-final --fields=sensor:list \
       | sed -e "s/,$SEN,/,/" \
25     | sed -e "s/^$SEN,/" | sed -e "s/,$SEN$//")

28     rfilter --sensor=$OTHERS $OUTBOUND --pass=stdout \
       | rwsset --sip-file=${SEN}-other.set --copy=stdout \
       | rwbag --bag-file=sipv4,flows,${SEN}-oth.bag

31     rwssettool --difference ${SEN}-out.set ${SEN}-other.set \
       --output=${SEN}-only.set
       rwbagtool --compare=ge ${SEN}-out.bag ${SEN}-oth.bag \
34     --output=${SEN}-most.bag
       rwbagtool --coverset ${SEN}-most.bag --output=${SEN}-most.set

37     rwssetcat --net=27 ${SEN}-most.set ${SEN}-only.set \
       | sed -e 's/|.*//' | rwssetbuild stdin cand.set

40     rfilter --sensor=$SEN $INBOUND --pass=stdout \
       | rwsset --dip-file=${SEN}-in.set

43     rwssettool --difference ${SEN}-in.set ${SEN}-most.set \
       ${SEN}-only.set --output=${SEN}-close.set
       rwssettool --union ${SEN}-only.set ${SEN}-most.set \
46     ${SEN}-close.set --output=${SEN}.set
       rm -f $FILES

49     echo "Sensor: $SEN -- End"
done

```

## Chapter 6

# Advanced Exploratory Analysis with SiLK: Exploring and Hunting

This chapter introduces advanced exploratory analysis through application of the analytic development process with the SiLK tool suite. It discusses the exploratory analysis process, starting points for exploration, and advanced SiLK commands and techniques that can be employed during an analysis.

Upon completion of this chapter you will be able to

- describe advanced exploratory analysis and how it maps to the analytic development process
- describe SiLK tools that often are used during exploratory analysis
- provide example workflows for exploratory network flow analysis

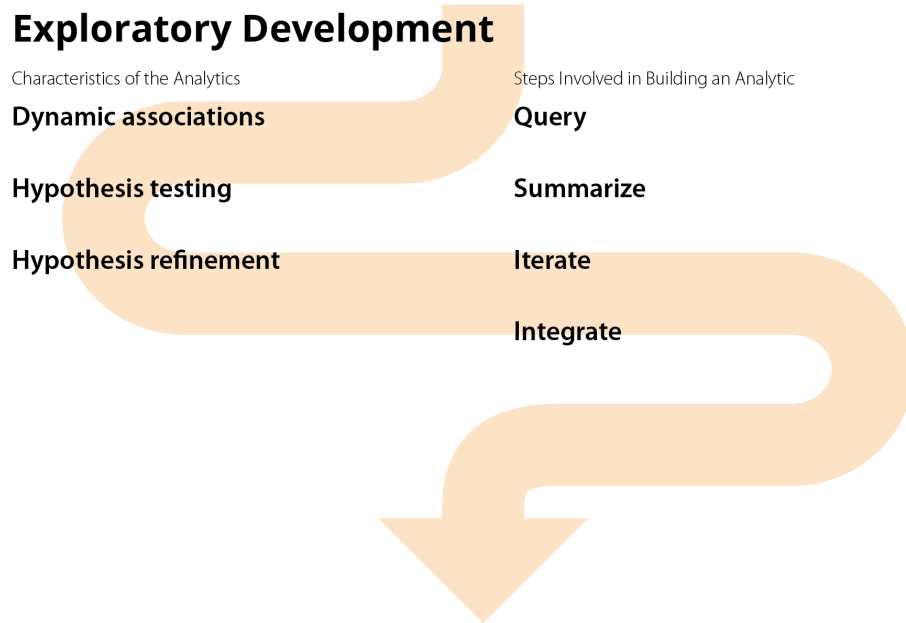
### 6.1 Exploratory Analysis: Concepts

The *exploratory* approach is the most open-ended of the analysis workflows described in this guide. It provides a framework for investigating more complex questions about network traffic. As the name suggests, exploratory analysis involves asking questions about trends, processes, and dynamic behavior that do not necessarily have fixed or obvious answers. Often the analysis leads to more questions!

Exploratory analysis uses single-path and multi-path analyses as building blocks to provide insight into network events. It typically considers more than one indicator of network behavior to provide a more complete understanding of what happened. Each building block represents a question (or part of a question) whose answer feeds into subsequent steps in the analysis. Analysts can assemble these building blocks by hand to prototype an analysis or examine one-time phenomena, iterating if necessary. This manual analysis can then transition to scripted analysis to save time, make it easier to repeat the analysis, and ensure more consistent results.

As Figure 6.1 indicates, this form of analysis rarely proceeds directly from initial question to final result. Instead, it gathers a variety of contextual information, and goes forward in a search pattern to reach the result.

Figure 6.1: Exploratory Analysis



### 6.1.1 Exploring Network Behavior

Exploratory analysis follows the general SiLK analysis workflow described in Figure 1.3. The overall course of exploratory analysis generally follows these steps:

1. come up with initial questions as the starting point (the Formulate stage)
2. apply a set of analyses to provide insight into a given question (the Model and Test stages)
3. integrate the results of these analyses with previously-available data (the Analyze stage)
4. refine the results (the Refine stage)
5. identify new questions to be investigated with further analyses (return to the Formulate stage to begin another iteration of the exploratory analysis)

In contrast to single-path and multi-path analysis, analysts explore with iterations of questions and their associated analyses. We may not have a specific set of behaviors, known event, or identified service in mind at the start of an exploratory analysis. Instead, we start with a general question and seek a clearer target during the iterations of the analysis.

As the exploration progresses, the scope of these questions often becomes more associated with any features of interest that are identified in the data. These features can be identified during the initial stages of establishing context for the network behavior to be investigated and gathering data about it. They also can emerge from the results of preliminary single-path or multi-path analyses of the data, which help to refine the scope of the exploratory analysis.

Exploratory analysis often takes a deep dive into the data to examine the possible causes or conditions related to the features of interest. The analyst may form, test, and either continue to investigate or abandon hypotheses about sources of the behavior. This exploration continues until the analyst reaches a sufficient level of confidence in understanding the features to either identify the causes of the event and share the results, or close the hypothesis. When all threatening hypotheses have been closed, we can close the investigation and associate the behavior with unremarkable network activity.

### 6.1.2 Starting Points for Exploratory Analysis

The first question that normally arises during an exploratory analysis is, “Where should it start?”

One set of starting points for exploratory analyses involves investigating unexplained activity on the monitored network. This produces initial questions such as

- What activity is present that could be malicious or damaging?
- What is the impact of this unexpected activity?
- What led to this unexpected activity?

The initial stage of the investigation may be a relatively straightforward profiling analysis or an inventory of traffic of interest across specific servers. Further iterations use the results of these profiles and inventories to identify outliers, inconsistencies, and behaviors exhibited by the unexplained activity. These outliers, inconsistencies, and behaviors then become the focus of exploration for the next iteration. At this point in the analysis, a case-by-case frequency summary or time-based trends can provide insight into the activity. The iterations of analysis continue until the analyst identifies a reasonable explanation of the activity, rejecting alternate explanations.

Another set of starting points for exploratory analyses involves unusual levels or directions of network services. This produces initial questions such as

- What led to these services?
- Is there a change in the population of external hosts employing these services?
- Does this change in service match with a model of network attack?

Analyses for these questions often summarize network activities via time series or service-by-service counts of incoming activity versus outgoing activity. The analyst integrates and examines these summaries and counts, looking for points of significant change in trends or traffic levels. These points of change, and the hosts involved in those changes, become the focus of the next iteration of the exploratory analysis.

### 6.1.3 Example Exploratory Analysis: Investigating Anomalous NTP Activity

As an example of an exploratory analysis, consider an investigation into Network Time Protocol (NTP) usage. NTP is associated with UDP port 123, so our starting point would be to find out what level of activity is observed on this port.

### Compare Inbound and Outbound UDP Activity on Port 123

The `rwfilter` commands in Example 6.1 provide an initial view of UDP activity on port 123. Command 1 profiles the outgoing activity to UDP/123. Command 2 profiles the incoming activity from this port. The size of the outgoing activity is surprisingly large (almost double the flow records, and one third larger in bytes) and worth looking at more closely.

Example 6.1: Using `rwfilter` to Profile NTP Activity

```
<1>$ rwfilter --start=2015/06/16 --sensor=S0,S1,S2,S3,S4 \
--type=out --proto=17 --dport=123 --pass=NTP-out.raw \
--print-vol=stdout
|          Recs |          Packets |          Bytes |          Files |
Total |      5566807 |      15791472 |      2646290347 |           60 |
Pass |         7040 |         8410 |         659800 |           |
Fail |      5559767 |      15783062 |      2645630547 |           |
<2>$ rwfilter --start=2015/06/16 --sensor=S0,S1,S2,S3,S4 \
--type=in --proto=17 --sport=123 --pass=NTP-in.raw \
--print-vol=stdout
|          Recs |          Packets |          Bytes |          Files |
Total |      3549008 |      15506147 |      3118288427 |           60 |
Pass |         3936 |         6148 |         493968 |           |
Fail |      3545072 |      15499999 |      3117794459 |           |
```

### Differentiate Benign and Malicious Traffic

The second phase of this sample exploratory analysis examines the UDP activity on port 123 more closely to look for patterns that might differentiate between benign and malicious traffic. `rwuniq` is well suited for this task, since it allows us to look at the traffic over a variety of contingencies. RFC 5905<sup>10</sup> specifies that NTP packets must have a byte size of 76 bytes for a request and 96 bytes for a response. We would therefore expect to find flow bytes that are a multiple of either of these two values.

Example 6.2 runs `rwuniq` on the output files from Example 6.1. Command 1 shows an excerpt of the various byte counts for outgoing flows. Command 2 shows a similar excerpt for the incoming flows. The byte sizes are the ones expected: 76, 96, 152 =  $2 \times 76$ , 192 =  $2 \times 96$ . (While not shown in the excerpt, all of the other byte sizes are multiples of either 76 or 96.)

What is surprising are the further columns on the last line of the excerpts in Example 6.2. NTP is a protocol where machines on a network make queries to a few servers for time values, then adjust their clocks based on the answers received. We would therefore expect outgoing traffic to come from only a few addresses but go to multiple hosts in the local network.

For flows with byte size 192 (double the response value), the number of local hosts and the number of remote hosts are nearly equal, which is unusual. Additionally, the number of incoming packets is larger than the number of outgoing requests. This disparity is also unusual, indicating that further exploration of this traffic is needed.

<sup>10</sup>Internet Engineering Task Force (IETF) Request for Comments 5905 (<https://www.ietf.org/rfc/rfc5905.txt>)

Example 6.2: Using `rwuniq` to examine NTP Activity

```

<1>$ rwuniq --fields=bytes \
  --values=Flows,packets,distinct:dip,distinct:sip --sort \
  NTP-out.raw \
| head -5
  bytes|   Records|   Packets|dIP-Distin|sIP-Distin|
   76|   5728|   5728|      2|      17|
   96|   362|   362|      6|      75|
  152|   563|  1126|      1|       1|
  192|   323|   646|     50|     46|
<2>$ rwuniq --fields=bytes \
  --values=Flows,packets,distinct:dip,distinct:sip --sort \
  NTP-in.raw \
| head -5
  bytes|   Records|   Packets|dIP-Distin|sIP-Distin|
   76|   2083|   2083|     16|       1|
   96|    66|    66|     27|       5|
  152|  1126|  2252|      1|       1|
  192|   615|  1230|     50|     46|

```

### Plot Anomalous Flows

Now that we have identified the anomalous flows, it's time to take a more detailed look at them. Example 6.3 uses `rwfilter` to isolate the incoming 192-byte flows, then calls `rwcount` to generate a time series of that traffic. The results are somewhat lengthy, since the `rwcount` command splits a day of traffic in five minute bins.

Command 1 uses `rwcount` options to generate a comma-separated-value file (CSV), which can be fed into Microsoft Excel or another graphing program. The resulting time series plot is shown in Figure 6.2.

NTP traffic is normally expected to either cluster around the point at which computers are connecting to the network, or occur periodically thereafter as the computers adjust their clocks. As can be seen in the figure, while there are regular pulses of traffic (visible more to the right end of the timeline), there are also irregular collections of traffic earlier in the day.

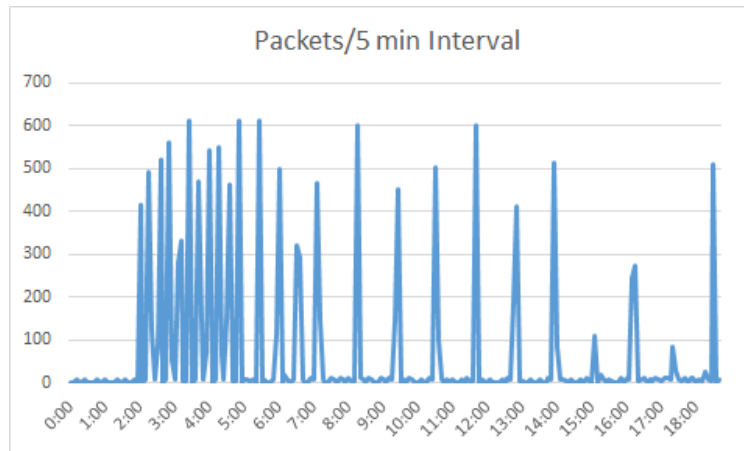
Example 6.3: Using `rwcount` to generate NTP Timelines

```

<1>$ rwfilter NTP-in.raw --bytes=192-192 --pass=stdout \
| rwcount --bin-size=300 --delim=',' >NTP-in.csv

```

Figure 6.2: Time Series Plot of NTP Traffic



### Is the Irregular NTP Traffic Related to Client-side Initialization?

To determine if this early irregular traffic is related to initialization of the client side of the NTP traffic, Example 6.4 generates summary bags for the hosts involved. Command 1 starts this process by identifying the earliest start time for each client (internal recipient) of the NTP traffic, using a combination of `rwfilter` and `rwuniq`. The output of the `rwuniq` is formatted for easy parsing in command 3, then saved to the file `source-fields.txt`.

An obvious approach at this point would be to issue another call to `rwfilter` to pull traffic prior to the start of the NTP traffic for the destinations involved, then call `rwuniq` to summarize the traffic. One complication to this approach is that the start times for the NTP traffic are not aligned. Rather than setting a fixed time (such as the minimum start time), it seems wise to analyze client traffic on a one-by-one basis. This individual approach has two disadvantages: it requires more programming (14 lines of commands, as opposed to two or three), and it is somewhat slower in execution (for the data being examined, about five seconds, rather than an immediate result). However, these disadvantages are outweighed by the increased precision of the results.

Command 2 of Example 6.4 prepares to analyze traffic for individual clients by creating an empty bag.

Command 3 is the loop that reads the start time information for each client (shown by the redirection at the end and the `read` call in the `while` statement's condition). The body of the loop then uses `date` as the NTP start time for the client and calculates a ten minute time window prior to the start time (using shell arithmetic and `awk` to reformat the range elements as SiLK formatted date-time values).

The bag initialized in command 2 (`cur.bag`) is updated for each source IP address. The loop in command 3 calls `rwfilter` using the time range for each address to find activity prior to the NTP traffic, using `rwbagtool` to add the results to the bag. Since `cur.bag` is initially empty but keyed by the client IP address, the prior activity is calculated for each client in turn. Although the `rwbagtool` call adds the bags, it is really just inserting the new entry in the summary bag, using `temp.bag` as an intermediate that is renamed in the `mv` call to overwrite the summary bag.

After all of the clients are summarized, command 4 displays the first few entries in `cur.bag`. They represent a brief sample of the 48 lines of full output.



---

 Example 6.4: Using `rwuniq` and Bags to Summarize Prior Traffic on NTP Clients
 

---

```

<1>$ rwfilter NTP-in.raw --bytes=192-192 --pass=stdout \
| rwuniq --fields=dip --values=stime-earliest --no-titles \
  --delim=' ' --output=source-fields.txt
<2>$ rwbagbuild --bag-input=/dev/null --key-type=sipv4 \
  --counter-type=records --output=cur.bag
<3>$ while IFS="" read -r line || [[ -n "$line" ]]; \
do srcArray=($line); \
  StEpoch=$( date -d ${srcArray[1]} +"%s"); \
  StTimeE=$(echo $(( $StEpoch - 1 )) \
| awk '{print strftime("%Y/%m/%dT%T",$1)}'); \
  StTimeS=$(echo $(( $StEpoch - 600 )) \
| awk '{print strftime("%Y/%m/%dT%T",$1)}'); \
  rwfilter --start=2015/06/16 --sensor=S0,S1,S2,S3,S4 \
  --type=out,outweb --saddress=${srcArray[0]} \
  --stime=${StTimeS}-${StTimeE} --pass=stdout \
  | rwbag --bag-file=sipv4,flows,stdout \
  | rwbagtool --add cur.bag stdin --output=temp.bag; \
  mv temp.bag cur.bag; \
done < source-fields.txt
<4>$ rwbagcat cur.bag | head -10
  10.0.40.54|          1452|
  192.168.40.25|       1777|
  192.168.40.50|        183|
  192.168.40.51|        123|
  192.168.50.11|         63|
192.168.111.109|       211|
192.168.111.131|       319|
  192.168.121.57|       325|
  192.168.121.77|       291|
192.168.121.145|       248|

```

---

- Low values (in the tens) would support the hypothesis that the host in question is in the process of initialization, as is possible for 192.168.50.11.
- High values (in the hundreds to thousands) would support that the host in question is in ongoing use, as is the case for most of the hosts.

### Possible Explanations for Anomalous NTP Traffic

While it is clear that there is a range of flow counts during the ten-minute time window prior to the start of NTP traffic, the numbers are large enough to eliminate the idea that the anomalous NTP traffic is due to initialization of the clients. The irregular traffic could be due to a variety of causes:

- It could be a covert means of mapping clients: send NTP results to the clients, then see which ones respond with an ICMP `host unreachable` message.
- It could be a covert signaling or command interface for malicious software that is running on the clients.
- It could reflect some flaw in the NTP implementation on the network. However, the relative maturity of NTP and its implementations and the halting of the irregular traffic discount this possibility.
- It could be a reflection attack against the target network. However, the nature of the exercise under way and the overall level of traffic involved discount that possibility.

Testing and evaluating these possibilities is left for further rounds of exploratory analysis!

### 6.1.4 Observations on Exploratory Analysis

The initial examination of inbound and outbound data in an exploratory analysis is the basis for all that follows. It is the starting point for an accumulating body of experience that aids the analyst in interpreting the results.

The odd behavior shown in Example 6.2 reflects a general aspect of exploratory analysis: it is shaped by the dynamics of network behavior. Since the analyst is exploring previously-unexamined behaviors, the process incorporates consideration of new behaviors and trends. This is done by both noting results that are expected in a given situation and those that are surprising.

Surprising results (such as those in Example 6.2) often serve as starting points for further rounds of analysis. Networks are not static: the traffic they carry constantly changes over time. Analysts therefore must build an evolving understanding of their behavior.

The exploratory analysis in this section uses a variety of SiLK tools. It also employs tabular data, graphical summary, and calculation of specific values. This diverse set of tools and techniques reflects the open-ended nature of exploratory analysis. Analysts need to obtain different views into the data to obtain a good understanding of the target of the analysis.

The outcome of an exploratory analysis can take several forms. As the analysis proceeds, it produces more and more background information on the network hosts and their traffic. As this information is captured, it aids in interpreting both results obtained later in the analysis and the results of other analyses.

Background information is not the only thing that can be reused in an exploratory analysis! Its component parts can be applied in other analyses as filters or analytics. For example, the date math portion of Example 6.4 is also incorporated into the case study discussed in Chapter 7. As the body of such filters and

analytics expands, you can perform further exploratory analyses more easily—simply reuse those developed for earlier analyses.

## 6.2 Exploratory Analysis: Analytics

The SiLK commands described in this chapter involve sophisticated uses of the SiLK tool suite that are often appropriate for exploratory analyses. However, they can be used with any analysis method.

### 6.2.1 Using Tuple Files for Complex Filtering

Partitioning criteria for many analyses comprise specific combinations of field values, any one of which can be considered as *passing*. While you can make repeated `rwfilter` calls and merge them later, this approach is often inefficient as it may involve pulling the same records from the repository several times.

For example, consider an analysis that is looking for a Simple Network Management Protocol (SNMP) call generated from viewing a malicious email message. SNMP is associated with UDP port 161, email receipt with TCP port 25. The naive approach would be a call to `rwfilter` that simply merges these two ports:

```
rwfilter --protocol=6,17 --dport=25,161
```

This call to `rwfilter` actually uses four permutations of the selection parameters to select records (protocol 6 and dport 25, protocol 6 and dport 161, protocol 17 and dport 25, protocol 17 and dport 161), not just the two that apply for this example (protocol 6 and dport 25, protocol 17 and dport 161).

As shown in Example 6.5, you could use two calls to `rwfilter` to pull the desired port-protocol permutations separately, then combine the data files with the `rwappend` command. However, this approach involves two calls to `rwfilter` that re-read the input flow records. If the analysis includes many cases, the same data would be read many times. On top of that, if the data set is large, each pull from the repository could take a significant amount of time.

Example 6.5: Using Multiple Data Pulls to Filter on Multiple Criteria

---

```
<1>$ rwfilter --start=2015/06/17 --type=in,out --protocol=6 \
  --dport=25 --pass=result.raw
<2>$ rwfilter --start=2015/06/17 --type=in,out --protocol=17 \
  --dport=161 --pass=part2.raw
<3>$ rwappend result.raw part2.raw
<4>$ rm part2.raw
```

---

A more efficient approach is to store partitioning criteria as a *tuple file* and use that file with `rwfilter` to pull the records in a single operation. A tuple file is a text file consisting of the five-tuple fields (`sIP`, `dIP`, `sPort`, `dPort`, `protocol`) delimited by vertical bars (`|`). `rwfilter` pulls the flow records that match the entries in the tuple file from the SiLK repository.

To select the protocol-dport combinations of (6,25) and (17,161), you could create a tuple file that contains both combinations:

```

protocol|dport
 6|25
17|161

```

Running `rwfilter` with the `--tuple-file` switch set to the name of this tuple file will select only those flow records with (protocol 6, dport 25) and (protocol 17, dport 161). Only one call to `rwfilter` would be needed to pull these records, not two (as in Example 6.5).

### Using Tuple Files to filter Web Servers

Example 6.6 shows a tuple file that is used to choose web server addresses on different ports.

- Command 1 shows the tuple file, `webservers.tuple`. The first line contains headers that identify the fields associated with the columns, `dIP` and `dPort`. The rest of the tuple file lists the destination IP address and destination port combinations of interest.
- This file can then be used with `rwfilter` as shown in command 2. The `--tuple-file` option need not be the only partitioning option. In command 2, the `--protocol` parameter also is specified as a partitioning criterion.

In some cases, you can obtain results more quickly by using seemingly redundant command-line options to duplicate some of the values from the tuple file. For instance, adding `--dport=80,443,8443` to the `rwfilter` call in Example 6.6 reduces the number of records that need to be examined with the tuple file. No matter where they appear in the `rwfilter` call, tuple files are always processed after the parameters that partition based on individual flow record fields, and before those using plug-ins or Python. However, filtering by multiple destination ports is not a substitute for the tuple file in Example 6.6, as these criteria are not sufficiently restrictive to produce the desired results.

## 6.2.2 Manipulating Bags

Using bags to store key values and counts (as described in Section 4.2.4) can be very helpful to store the intermediate and final results of a multi-path analysis. SiLK supports several advanced options for working with bags. In addition to comparing bags with `rwbagtool` (Section 4.2.4) and extracting sets from bags (Section 4.2.5), you can use `rwbagtool` to do the following:

- add and subtract bags (analogous to the SiLK set operations)
- multiply bags by scalar numbers
- divide bags
- threshold bags (filter bags on volume)

The result of these operations is a bag with new volumes.

Example 6.6: Filtering on Multiple Criteria with a Tuple File

```

<1>$ cat <<EOF >webservers.tuple
      dIP|dPort
    10.0.40.21| 443
    10.0.40.23| 8443
    10.0.20.59|  80
  192.168.20.59|  80
    10.0.40.21|  80
  192.168.40.24| 443
  192.168.40.27| 443
  192.168.40.91| 443
  192.168.40.92| 443
EOF
<2>$ rfilter --type=in,inweb --start-date=2015/06/02 \
  --end-date=2015/06/18 --dport=80,443,8443 \
  --protocol=6 --flags-all=SAF/SAF,SAR/SAR \
  --tuple-file=webservers.tuple --sensors=S0,S1 \
  --pass=stdout \
| rwuniq --fields=dIP,dPort --value=Records
      dIP|dPort|  Records|
    10.0.20.59|  80|    29720|
    10.0.40.21| 443|   355791|
    10.0.40.23| 8443|   30934|
  192.168.40.91| 443|      6|
  192.168.40.27| 443|      9|
  192.168.40.92| 443|      3|
  192.168.20.59|  80|   10652|
    10.0.40.21|  80|    1565|
  192.168.40.24| 443|     24|

```

## Adding and Subtracting Bags

Suppose you want to find the total number of records associated with the IP addresses that are stored in two bags. You can add the contents of the two bags together to create a new bag that holds the sum of their contents.

To add bags together, use the `rwbagtool --add` parameter. The `--output-path` parameter specifies where to deposit the results. Most of the results from `rwbagtool` are bags themselves. Example 6.7 shows how to use bag addition to find the total number of flows of inbound web traffic over a two-day period.

1. The `rwbagcat` calls in commands 1 and 2 display the contents of the two bags to be added, `web-20150616.bag` and `web-20150617.bag`. Each bag contains a day's worth of inbound web traffic flows.
2. Command 3 adds the two bags with the `rwbagtool --add` parameter.
3. The results of the addition are stored in `web-sum.bag`, which is shown in command 4.
  - If an IP address appears in both bags, the `rwbagtool --add` command sums up the number of flows in the two bags. For instance, the IP address `10.0.20.59` appears in both bags. The number of flows for this IP address in `web-sum.bag` is the sum of the flows in the two bags.
  - If an IP address appears in just one bag, the `rwbagtool --add` command still includes it in the results. For instance, the IP address `190.168.40.27` only appears in the bag `web-20150616.bag` but is included in the results stored in `web-sum.bag`.

Similarly, you may want to subtract the byte counts in one bag from those stored in another bag to find out how many are left over after a step in your analysis. Use the `rwbagtool --subtract` command to subtract the contents of bags.

Bag subtraction operates in the same fashion as bag addition, but all bags after the first are subtracted from the first bag specified in the command. Bags cannot contain negative values: any subtraction resulting in a negative number causes `rwbagtool` to omit the corresponding key from the resulting bag.

**Hint:** Bags do store information in the file header about which types of keys and counts they contain. However, the information is not used to restrict bag operations. Consequently, `rwbagtool` will add or subtract byte bags and packet bags without warning, producing meaningless results.

If unequal but compatible types are added or subtracted, a meaningful result type will be produced. For example, keys of `sIPv4` and `dIPv4` will produce a result key of type `any-IPv4`. When incompatible types are combined, the resulting type will be `custom` (the most generic bag type). Use `rwfileinfo --fields=bag` to view this information, as described in Section 4.2.9.

## Multiplying and Dividing Bags

You can multiply the values in a bag by a scalar number and divide the contents of a bag by the contents of another bag. This is useful for operations such as computing percentages (for instance, to compare traffic levels during different time periods).

- Use the `rwbagtool --scalar-multiply` command to multiply all of the counter values in a bag by a scalar value. Bags can only be multiplied by a scalar value.

Example 6.7: Merging the Contents of Bags Using `rwbagtool --add`

---

```
<1>$ rwbagcat web-20150616.bag
  10.0.20.59|          7977|
  10.0.40.21|        135757|
  10.0.40.23|        11700|
  192.168.20.59|       3980|
  192.168.40.24|         17|
  192.168.40.27|         9|
  192.168.40.91|         3|
<2>$ rwbagcat web-20150617.bag
  10.0.20.59|        15248|
  10.0.40.21|       221599|
  10.0.40.23|       19234|
  192.168.20.59|       6672|
  192.168.40.24|         7|
  192.168.40.91|         3|
  192.168.40.92|         3|
<3>$ rwbagtool --add web-20150616.bag web-20150617.bag \
>web-sum.bag
<4>$ rwbagcat web-sum.bag
  10.0.20.59|        23225|
  10.0.40.21|       357356|
  10.0.40.23|       30934|
  192.168.20.59|      10652|
  192.168.40.24|         24|
  192.168.40.27|         9|
  192.168.40.91|         6|
  192.168.40.92|         3|
```

---

- Use the `rwbagtool --divide` command to divide the counter values in one bag by those of another.

**Hint:** Be very careful when dividing bags. **The second (denominator) bag must contain every key found in the first (numerator) bag—do not divide by zero!** To ensure that the elements of the two bags match, use the `rwbagtool --intersect` command to remove mismatched elements.

1. Extract the set of IP addresses from the denominator bag by using `rwbagtool --coverset` as described in Section 4.2.5.
2. Run `rwbagtool --intersect` on the numerator bag to remove all elements that are not found in the denominator bag as described in Section 4.2.5.
3. Use `rwbagtool --divide` to divide the contents of the numerator bag by those of the denominator bag.

Example 6.8 shows how to use scalar multiplication and division. The example computes the percentage change in traffic between the two bags from Example 6.7. It uses `rwbagtool --coverset` to remove the IP addresses that do not appear in both bags, then uses the `rwbagtool` options `--scalar-multiply` and `--divide` to compute the percentage change in traffic between the IP addresses in both bags.

### Thresholding Bags with Count and Key Parameters

Sometimes, you may want to threshold the contents of a bag to items that are larger than a minimum value or smaller than a maximum value. This thresholding can be used to limit key values (e.g., eliminating IP addresses that are lower or higher than the specified address value) as well as count values (e.g., eliminating IP addresses with packet counts that are lower than the desired volume).

The `--minkey`, `--maxkey`, `--mincounter`, and `--maxcounter` parameters supported by `rwbagcat` are also supported by `rwbagtool`. In this case, they specify the minimum and maximum key and count values for output. They can optionally be combined with an operation parameter (e.g., `--add`, `--subtract`) or a masking parameter (i.e., `--intersect` or `--complement-intersect`) to perform other operations on a bag. Example 6.10 shows an example of thresholding by minimum and maximum counts.

### 6.2.3 Sets Versus Bags: A Scanning Example

Both sets and bags can be employed to search for network scanners. This section provides some examples of each and contrasts how they are used within an analysis.

#### Fine-tuning IP Sets to Find Scanners

Using IP sets can focus on alternative representations of traffic and identify network scanning and other activities. Example 6.9 drills down on IP sets themselves and provides a different view of this traffic.

This example isolates the set of hosts that exclusively scan from a group of flow records using `rwfilter` to separate the set of IP addresses that complete legitimate TCP sessions from the set of IP addresses that never complete sessions. As this example shows, the `final.set` set file consists of two IP addresses in contrast to the set of thirty-six that produced low-packet flow records—these addresses are consequently suspicious.<sup>11</sup>

<sup>11</sup>While this might be indicative of scanning activity, the task of scan detection is more complex than shown in Example 6.9.



Example 6.8: Using `rwbagtool` to Generate Percentages

```

<1>$ rwbagtool --coverset 20150616.bag >20150616.set
<2>$ rwbagtool --coverset 20150617.bag >20150617.set
<3>$ rwsettool --intersect 20150616.set 20150617.set \
    >common.set
<4>$ rwbagtool --scalar-multiply=100 --intersect=common.set \
    20150616.bag >multiply.bag
<5>$ rwbagcat multiply.bag
    10.0.20.59|          797700|
    10.0.40.21|         13575700|
    10.0.40.23|         1170000|
    192.168.20.59|       398000|
    192.168.40.24|        1700|
    192.168.40.91|        300|
<6>$ rwbagcat 20150617.bag
    10.0.20.59|          15248|
    10.0.40.21|         221599|
    10.0.40.23|         19234|
    192.168.20.59|       6672|
    192.168.40.24|         7|
    192.168.40.91|         3|
    192.168.40.92|         3|
<7>$ rwbagtool --intersect=common.set 20150617.bag \
    >predivide.bag
<8>$ rwbagtool --divide multiply.bag predivide.bag >divide.bag
<9>$ rwbagcat divide.bag
    10.0.20.59|          52|
    10.0.40.21|          61|
    10.0.40.23|          61|
    192.168.20.59|        60|
    192.168.40.24|       243|
    192.168.40.91|       100|

```

Example 6.9: Using `rwset` to Filter for a Set of Scanners

```

<1>$ rfilter --start-date=2015/06/02 --protocol=6 \
    --type=in,inweb --packets=1-3 --pass=stdout \
| rwset --sip-file=low.set
<2>$ rfilter --start-date=2015/06/02 --protocol=6 \
    --type=in,inweb --packets=4- --pass=stdout \
| rwset --sip-file=high.set
<3>$ rwsettool --difference low.set high.set \
    --output-path=final.set
<4>$ rwsetcat low.set --count-ips
36
<5>$ rwsetcat final.set --count-ips
2

```

## Using Bags to Find Scanners

To show how bags differ from sets, let's revisit the scanning filter presented in Example 6.9. The difficulty with that code is that if a scanner completed *any* handshake, it would be excluded from the `low.set` file. Many automated scanners would fall under this exclusion if any of their potential victims responded to the scan. It would be more robust to include as scanners hosts that complete only a small number of their connections (10 or fewer) and have a reasonable number of flow records covering incomplete connections (10 or more).

By using bags, Example 6.10 is able to incorporate counts, resulting in the detection of more potential scanners.

1. The calls to `rwfilter` in commands 1 through 3 are piped to `rwbag` to create bags for incomplete, FIN-terminated, and RST-terminated traffic.
2. Commands 4 and 5 use `rwbagtool --coverset` to generate the cover sets for these bags. These commands also use thresholding to generate two sets: `fast-low.set`, which contains only IP addresses with fewer than 10 low-packet connections (`--mincounter=10`) and `fast-high.set`, which contains only IP addresses with more than 10 incomplete connections (`--maxcounter=10`).
3. Command 6 uses `rwsettool --difference` to find the set of IP addresses that are members of `fast-low.set` but not members of `fast-high.set`. The result, `scan.set`, represents the set of IP addresses that responded to the scans.
4. Command 7 uses `rwsetcat` to count the number of IP addresses in each bag.

## 6.2.4 Manipulating SiLK Files

### Combining Flow Record Files to Provide Context

When you are profiling flow records, you may want to drill down into the data to find specific behaviors. This is especially useful for analyzing traffic with large volumes (for example, by the duration of transfer and by protocol). Issuing repeated `rwfilter` calls subdivides large data sets into smaller ones that are easier to examine and manipulate. However, sometimes this obscures the “big picture” of what is happening during an event. Combining flow record files is one way to provide this kind of context.

Use one of the following SiLK commands to merge multiple flow record files:

- `rwcat` concatenates flow record files in the order in which they are listed. It creates a new flow record file that contains the merged records.
- `rwappend` places the contents of the flow record files at the end of the first specified flow record file. It does not create a new file.

---

Scanners sometimes complete connections to hosts that respond (to exploit vulnerable machines); non-scanning hosts sometimes consistently fail to complete connections to a given host (contacting a host that no longer offers a service). A more complete set of scan detection heuristics is implemented in the `rwscan` tool, which is discussed in Section 4.2.4.

Example 6.10: Using `rwbagtool` to Filter Out a Set of Scanners

---

```

<1>$ rfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --type=in,inweb --bytes=2048- --pass=stdout \
| rfilter stdin --duration=1200- --pass=slowfile.rw \
  --fail=fastfile.rw
<2>$ rfilter fastfile.rw --protocol=6 --flags-all=S/SRF \
  --packets=1-3 --pass=stdout \
| rwbag --sip-flows=fast-low.bag
<3>$ rfilter fastfile.rw --protocol=6 \
  --flags-all=SAF/SARF,SR/SRF --pass=stdout \
| rwbag --sip-flows=fast-high.bag
<4>$ rwbagtool fast-low.bag --mincounter=10 --coverset \
  --output-path=fast-low.set
<5>$ rwbagtool fast-high.bag --maxcounter=10 --coverset \
  --output-path=fast-high.set
<6>$ rwssetool --difference fast-low.set fast-high.set \
  --output-path=scan.set
<7>$ rwssetcat fast-low.set fast-high.set scan.set --count-ips
fast-low.set:40
fast-high.set:104
scan.set:37

```

---

**Hint:** As an alternative to combining flow files, many of the SiLK tools accept multiple flow record files as input to a single call. For example, `rfilter` can accept several flow files to filter during a single call and `rwsort` can accept several flow files to merge and sort during a single call. Very often, it is more convenient to use multiple inputs than to combine flow files.

In Example 6.11, `rwcat` is used to combine previously filtered flow record files to permit the counting of overall values.

1. The initial call to `rfilter` in command 1 pulls out all records in the period of interest: 2015/6/10 through 2015/6/24. Subsequent calls to `rfilter` split these records into three files, depending on the duration of the flow:
  - slow flows of at least 20 minutes duration (i.e., those that match the partitioning criteria of `--duration=1200-`)
  - medium flows of 1–20 minutes duration (i.e., those that match the partitioning criteria of `--duration=60-1199`),
  - fast flows of less than 1 minute duration (the remainder of the flows, which had failed both partitioning criteria and must therefore be less than one minute long)
2. The calls to `rfilter` in Commands 2 through 4 split each of the initial divisions based on protocol: UDP (17), TCP (6), and ICMPv4 (1). They are saved to files whose names correspond to their speed and protocol (e.g., `slow17.rw`, `med17.rw`, `fast17.rw`).
3. The calls to `rwcat` in commands 5–7 combine the three splits for each protocol into one overall file per protocol (e.g., `all17.rw`).

4. Command 8 is a short script that produces a summary output reflecting the volume of records in each of the composite files.

When using the `rwfileinfo` command, be aware that `rwcat` creates a new file and can record annotation (using `--note-add` and `--note-file-add`) in the output file header. However, it does not preserve this information from its input files. `rwappend` cannot add annotations and command history to the output file.

For more information about the `rwcat` command, see Appendix C.10 or enter the command `rwcat --help`.

For more information about the `rwappend` command, see Appendix C.11 or enter the command `rwappend --help`.

### 6.2.5 Dividing or Sampling Flow Record Files with `rwsplit`

In addition to being able to join flow record files, some analyses are facilitated by dividing or sampling flow record files. To facilitate coarse parallelism, one approach is to divide a large flow record file into pieces and concurrently analyze each piece separately. For extremely high-volume problems, analyses on a series of robustly taken samples can produce a reasonable estimate using substantially fewer resources. `rwsplit` is a tool that facilitates both of these approaches to analysis.

Each call to `rwsplit` requires the `--basename` switch to specify the base file name for output files. In addition, one of these parameters must be present:

- `--ip-limit` specifies the IP address count at which to begin a new output file
- `--flow-limit` specifies the flow count at which to begin a new output file
- `--packet-limit` specifies the packet count at which to begin a new output file
- `--byte-limit` specifies the byte count at which to begin a new output file

Example 6.12 is an example of a coarsely parallelized process.

1. Command 1 pulls a large number of flow records with the `rwfilter` command, then use the `rwsplit` command to divide those records into a series of 400,000-record files with a base name of `part`.
2. Command 2 initializes a list of generated filenames. It then uses the `rwfilter` command to separate the records in each file based on sets that contain IP addresses of interest (`mission.set`, `threat.set`, `casual.set`).
3. In command 3, each of these files is then fed into an `rwfileinfo` call to count the number of records that fall into the selection categories (`mission`, `threat`, and `casual`).

Example 6.13 is an example of a sampled-flow process. These commands estimate the percentage of UDP traffic moving across a large infrastructure over a workday.

1. Command 1 invokes `rwfilter` to perform the initial data pull, retrieving a very large number of flow records. It then uses the `rwsplit` command to pull 100 samples of 1,000 flow records each (`--flow-limit=1000 --sample-ratio=100`), with a 1% rate of sample generation (that is, of 100 samples of 1,000 records, only one sample is retained).

Example 6.11: Combining Flow Record Files with `rwcat` to Count Overall Volumes

```

<1>$ rwwfilter --type=in,inweb --start-date=2015/6/10 \
  --end-date=2015/6/24 --protocol=0- --note-add='example' \
  --pass=stdout \
| rwwfilter stdin --duration=1200- --pass=slowfile.rw \
  --fail=stdout \
| rwwfilter stdin --duration=60-1199 --pass=medfile.rw \
  --fail=fastfile.rw
<2>$ rwwfilter slowfile.rw --protocol=17 --pass=slow17.rw \
  --fail=stdout \
| rwwfilter stdin --protocol=6 --pass=slow6.rw --fail=stdout \
| rwwfilter stdin --protocol=1 --pass=slow1.rw
<3>$ rwwfilter medfile.rw --protocol=17 --pass=med17.rw \
  --fail=stdout \
| rwwfilter stdin --protocol=6 --pass=med6.rw --fail=stdout \
| rwwfilter stdin --protocol=1 --pass=med1.rw
<4>$ rwwfilter fastfile.rw --protocol=17 --pass=fast17.rw \
  --fail=stdout \
| rwwfilter stdin --protocol=6 --pass=fast6.rw --fail=stdout \
| rwwfilter stdin --protocol=1 --pass=fast1.rw
<5>$ rwcat slow17.rw med17.rw fast17.rw --output-path=all17.rw
<6>$ rwcat slow1.rw med1.rw fast1.rw --output-path=all1.rw
<7>$ rwcat slow6.rw med6.rw fast6.rw --output-path=all6.rw
<8>$ echo -e "\nProtocol, all, fast, med, slow"; \
for p in 6 17 1; \
do rm -f ,c.txt ,t.txt ,m.txt
echo " count-records -" >,c.txt
for s in all fast med slow; \
do rwwfileinfo $$s$p.rw --fields=count-records \
| tail -n1 >,t.txt
join ,c.txt ,t.txt >,m.txt
mv ,m.txt ,c.txt
done
sed -e "s/^ *count-records - */$p,/" \
-e "s/\([^, ]*\),* */\1, /g" ,c.txt
done

Protocol, all, fast, med, slow
6, 6327373, 6280726, 46274, 373,
17, 7304579, 7258750, 45029, 800,
1, 131843, 124221, 6271, 1351,

```

Example 6.12: rwsplit for Coarse Parallel Execution

```

<1>$ rwsplit --type=inweb,outweb --start-date=2015/6/10 \
  --end-date=2015/6/24 --bytes-per-packet=45- \
  --pass=stdout \
| rwsplit --flow-limit=400000 --basename=part
# keep track of files generated
<2>$ s_list=(skip); \
for f in part*; \
do n=$(basename $f); \
t=${n%.*}; \
rm -f ${t{-miss,-threat,-casual,-other}.rw}; \
rwsplit $f --anyset=mission.set --pass=${t-miss.rw} \
  --fail=stdout \
| rwsplit stdin --anyset=threat.set --pass=${t-threat.rw} \
  --fail=stdout \
| rwsplit stdin --anyset=casual.set --pass=${t-casual.rw}; \
s_list=${s_list[*]} ${t{-miss,-threat,-casual}.rw}; \
done
<3>$ echo -e "\nPart-name, mission, threat, casual"; \
prev=" "; \
for f in ${s_list[*]}; \
do if [ "$f" = skip ]; \
then continue; \
fi; \
cur=${f%-*}; \
if [ "$prev" != " " ]; \
then if [ "$cur" != "$prev" ]; \
then echo; \
echo -n "$cur, "; \
fi; \
else echo -n "$cur, "; \
fi; \
prev=$cur; \
echo -n "$(rwsplit --fields=count-records $f | tail -n1 \
| sed -e "s/^ *count-records *//")", "; \
done; \
echo

```

```

Part-name, mission, threat, casual
part.00000000, 342291, 2191, 5741,
part.00000001, 337363, 2331, 6036,
part.00000002, 351400, 1365, 3403,
part.00000003, 324637, 1438, 4880,
part.00000004, 59355, 9590, 45748,
part.00000005, 32648, 9042, 43983,
part.00000006, 42775, 13484, 56446,
part.00000007, 52733, 11759, 62119,
part.00000008, 45304, 13089, 60488,
part.00000009, 61153, 3854, 19262,

```

2. Commands 2 through 4 create a file to store the summary results (`udpsample.txt`), then use the `rwstats` command to summarize each sample and isolate the percentage of UDP traffic (protocol 17) in the sample. The results in `udpsample.txt` are then sorted using the operating system `sort` command.
3. Commands 5 through 7 profile the resulting percentages to report the minimum, maximum, and median percentages of UDP traffic.

---

Example 6.13: `rwsplit` to Generate Statistics on Flow Record Files

---

```
<1>$ rfilter --type=in,inweb --start-date=2015/6/10 \
  --end-date=2015/6/24 --proto=0-255 --pass=stdout \
| rwsplit --flow-limit=1000 --sample-ratio=100 \
  --basename=sample --max-outputs=100
<2>$ echo -n >udpsample.txt
<3>$ for f in sample*rwf; \
do rwstats $f --values=records --fields=protocol --count=30 \
  --top \
| grep "17|" | cut -f3 "-d|" >>udpsample.txt
done
<4>$ sort -nr udpsample.txt >tmp.txt
<5>$ echo -n "Max UDP%: "; \
  head -n 1 tmp.txt
Max UDP%: 83.700000
<6>$ echo -n "Min UDP%: " ; \
  tail -n 1 tmp.txt
Min UDP%: 1.700000
<7>$ echo -n "Median UDP%: "; \
head -n 50 tmp.txt \
| tail -n 1
Median UDP%: 68.800000
```

---

For more information about the `rwsplit` command, see Appendix C.12 or enter the command `rwsplit --help`.

### 6.2.6 Generate Flow Records From Text

The `rtuc` (Text Utility Converter) tool creates SiLK flow record files from columnar text. `rtuc` is effectively the inverse of `rcut`, with additional parameters to supply values not given by the columnar input.

`rtuc` is useful when you need to work with tools or scripting languages that manipulate text output. For example, some scripting languages (Perl in particular) have string-processing functions that may be useful during an analysis. However, for later processing, you may need to use a binary file format for compactness and speed. In this situation, you could use `rcut` to convert the binary flow record files to text, process the resulting file with a scripting language, and then use `rtuc` to convert the text output back to the binary flow record format.

If scripting can be done in the Python programming language, the programming interface contained in the `silk` module allows direct manipulation of the binary flow records without converting them to text (and

back again). This binary manipulation is more efficient than text-based scripting.<sup>12</sup> See Chapter 8 for more information on using Python with SiLK.

On the other hand, `rwtuc` gives you complete control of the binary representation's content. This is very useful if you need to cleanse a flow record file before exchanging data<sup>13</sup> (for instance, to anonymize IP addresses). To ensure that unreleasable content is not present in the binary form, an analyst can convert binary flow records to text, perform any required edits on the text file, and then generate a binary representation from the edited text. Example 6.14 shows a sample use of `rwtuc` for anonymizing flow records. After `rwtuc` is invoked in command 3, both the file-header information and non-preserved fields have generic or null values.

Example 6.14: Simple File Anonymization with `rwtuc`

```
<1>$ rfilter --sensor=S0 --type=in --start-date=2015/06/02 \
--end-date=2015/06/18 --protocol=17 \
--bytes-per-packet=100- --pass=bigflows.rw
<2>$ rwcut bigflows.rw --fields=1-5,sTime --num-recs=20 \
| sed -re 's/([0-9]+\.)\{3\}/192.168.200./g' \
>anonymized.rw.txt
<3>$ rwtuc anonymized.rw.txt --output-path=anonymized.rw
<4>$ rfileinfo anonymized.rw
anonymized.rw:
  format(id)          FT_RWIPV6ROUTING(0x0c)
  version             16
  byte-order          littleEndian
  compression(id)     lzolx(2)
  header-length       88
  record-length       88
  record-version      1
  silk-version        3.16.0
  count-records       20
  file-size           512
  command-lines
      1 rwtuc --output-path=anonymized.rw anonymized.rw.txt
<5>$ rwcut anonymized.rw --fields=sIP,dIP,sTime,sensor \
--num-recs=4
  sIP|          dIP|          sTime|sen|
192.168.200.205| 192.168.200.20|2015/06/16T12:50:02.144| S0|
 192.168.200.5| 192.168.200.20|2015/06/16T12:50:03.139| S0|
192.168.200.218| 192.168.200.20|2015/06/16T12:50:05.189| S0|
192.168.200.160| 192.168.200.20|2015/06/16T12:50:09.997| S0|
```

`rwtuc` expects input in the default format for `rwcut` output. The record fields should be identified either in a heading line or in a `--fields` parameter of the call. `rwtuc` has a `--column-separator` parameter, with an argument that specifies the character-separating columns in the input. For debugging purposes, input lines that `rwtuc` cannot parse can be written to a file or pipe which the `--bad-input-lines` option names. For

<sup>12</sup>In several published examples, analysts encoded non-flow information as binary flow records using `rwtuc` or PySiLK so that SiLK commands could be used for the fast filtering and processing of that information.

<sup>13</sup>Ensuring that data content can be shared is quite complex, and involves many organization-specific requirements. `rwtuc` helps with mechanics, but often more transformations are required. The `rwtsetool` command contains parameters ending in `-strip` that also help to cleanse IP sets.



fields not specified in the input, an analyst can either let them default to zero (as shown in Example 6.14, especially for `sensor`) or use parameters of the form `--FixedValueParameter=FixedValue` to set a single fixed value for that field in all records, instead of using zero. Numeric field IDs are supported as arguments to the `--fields` parameter, not as headings in the input file.

For more information about the `rwttuc` command, see Appendix C.13 or enter the command `rwttuc --help`.

### 6.2.7 Labeling Data with Prefix Maps

Some analyses are easier to conceptualize and perform when you assign a text label to data of interest. You can identify this data and think of it in context by the label you've given it, not just as an abstract group of flow records. You can then filter and perform other operations on the data according to how it is labeled.

SiLK allows you to create a *prefix map* (often abbreviated as *pmap*) to assign user-defined text labels to ranges of IP addresses or protocols and ports. You can use the resulting pmap file to retrieve, partition, sort, count, and perform other operations on network flow records by their labels. This enables you to work with data semantically.

#### What Are Prefix Maps?

A prefix map file is a binary file that maps a value (either an IP address or a protocol-port pair) to a text label. SiLK supports two general types of prefix maps.

- **User-defined prefix maps** assign arbitrary text labels to IP addresses or protocol-port combinations.
- **Predefined prefix maps** are specialized prefix maps that are typically included with the SiLK distribution and facilitate analysis by country code or traffic direction.

Both types of pmap can be used interchangeably with the `rwfilter`, `rwcut`, `rwsort`, `rwuniq`, `rwstats`, `rwgroup`, and `rwmaplookup` commands to perform operations on SiLK network data according to how it is labeled.

**Comparing Prefix Maps to Sets and Bags.** Like SiLK IP sets and bags, prefix maps allow you to create user-defined groups of IP addresses to facilitate further analysis. However, prefix maps are more generalized groupings than sets and bags. Where a set creates a binary association between an IP address and a condition (an address is either in the set or not in the set), and a bag between an IP address and a numeric value, a prefix map assigns arbitrary, user-defined text labels to many different address ranges. Prefix maps also expand the type of groupings to include assigning labels to protocol-port ranges.

It is often easier to examine IP addresses by label rather than by whether they belong to a bag or set. For example, suppose you want to look at traffic from IP addresses that are linked to different types of malware. You could create multiple IP sets or bags that contain the addresses associated with each type of malware and compute traffic statistics individually for each set or bag. However, this would be cumbersome and time-consuming to script.

Alternatively, you could create a single, user-defined pmap file that labels the addresses by the type of malware that is associated with each address. You could then use the pmap file to filter network flow data and compute traffic statistics according to the type of malware, all in a single analytic. This is a more intuitive and easier way to perform that type of analysis.

## Creating a User-defined Prefix Map From a Text File

To create a binary prefix map from a text file, use the `rwpmmapbuild` tool. Creating a user-defined pmap is a two-step process:

1. Create a text file that contains the mapping of IP addresses or protocol-port pairs to their labels.
2. Use the `rwpmmapbuild` command to translate this text-based file into a binary pmap file.

**Creating the Text File.** Each mapping line in the text file contains either an IP address or protocol-port pair range with a corresponding label. They are separated by whitespace (spaces and tabs). Include the following information when creating a text file for use with the `rwpmmapbuild` command:

- The input file may specify a name for the pmap via the line `map-name mapname` where *mapname* is the name of the pmap. Pmap names cannot contain whitespaces, commas, or colons.
- `mode` specifies the type of pmap.
  - `ipv4` creates a pmap containing IPv4 addresses
  - `ipv6` creates a pmap containing IPv6 addresses
  - `proto-port` creates a pmap containing protocol-port pairs
- If you are creating an IP address pmap, specify an address range with either a CIDR block or a whitespace-separated low IP address and high IP address (formatted in canonical form or as integers). Specify a single host as a CIDR block with a prefix length of 32 for IPv4 or 128 for IPv6.
- If you are creating a protocol-port pmap, specify a range that is either a single protocol or a protocol and a port separated by a slash character (/). If the range is a single port, specify that port number as the starting and ending value of the range. For example, `17/17` specifies general UDP traffic; `17/53` specifies DNS traffic (UDP on port 53), and `17/67 17/68` specifies DHCP client and server traffic (UDP on ports 67 and 68).
- Do not use commas, which invalidate the pmap for use with `rwfilter`.
- Comment lines begin with the pound or hashtag character (#). Do not use this character in text labels.
- The input file may also contain a *default label* to be used when there is no matching range in the text file. This default is specified by the line `default deflabel`, where *deflabel* is the text label specified by the analyst for otherwise-unlabeled address ranges.

For more information about the `rwpmmapbuild` command, see Appendix [C.23](#) or enter the command `rwpmmapbuild --help`.

**Building the Prefix Map File.** Example 6.15 shows an example of how to create a prefix map of FCC network labels from the FCCX dataset described in Section 1.7. The FCC network description is contained in the file `fccnets.pmap.txt`. It associates network address ranges with text labels that identify their subnetwork locations (`Div0Ext`, `Div1Ext`, etc.).

In addition to the address list, the text file specifies the prefix map name (`map-name fccnets`). This name is used in SiLK commands to identify which prefix map is being used. Using the map name as the name of the text file and resulting pmap file helps you to keep track of and organize your user-defined prefix maps. Note also that the text file includes a default label (`default None`) that is assigned to IP addresses that are not listed in the FCC network description.

The `rwpmaphbuild` command takes `fccnets.pmap.txt` as input to create a binary pmap file, `fccnets.pmap`.

Example 6.15: Using `rwpmaphbuild` to Create a FCC Pmap File

---

```
<1>$ cat <<-EOF >fccnets.pmap.txt
map-name fccnets
default None

# FCC network descriptions
10.0.10.0/24 Div0Ext
10.0.20.0/24 Div0Ext
10.0.30.0/24 Div0Ext
10.0.40.0/24 Div0Ext
10.0.50.0/24 Div0Ext
192.168.10.0/24 Div1Ext
192.168.20.0/24 Div1Ext
192.168.30.0/24 Div1Ext
192.168.40.0/24 Div1Ext
192.168.50.0/24 Div1Ext
192.168.60.0/24 Div1Ext
192.168.70.0/24 Div1Ext
192.168.110.0/23 Div1Ext
192.168.120.0/23 Div1Ext
192.168.122.0/23 Div1Ext
192.168.124.0/24 Div1Ext
192.168.130.0/23 Div1Ext
192.168.140.0/23 Div1Ext
192.168.142.0/23 Div1Ext
192.168.150.0/23 Div1Ext
192.168.160.0/23 Div1Ext
192.168.162.0/23 Div1Ext
192.168.164.0/23 Div1Ext
192.168.166.0/24 Div1Ext
192.168.170.0/24 Div1Ext
10.0.40.0/24 Div0Int
10.0.50.0/24 Div0Int
192.168.20.0/24 Div1Int1
192.168.40.0/24 Div1Int1
192.168.50.0/24 Div1Int1
192.168.60.0/24 Div1Int1
192.168.70.0/24 Div1Int1
192.168.110.0/23 Div1Int1
192.168.120.0/23 Div1Int1
```

```
192.168.122.0/23 Div1Int1
192.168.124.0/24 Div1Int1
192.168.130.0/23 Div1Int1
192.168.140.0/23 Div1Int1
192.168.142.0/23 Div1Int1
192.168.150.0/23 Div1Int1
192.168.160.0/23 Div1Int1
192.168.162.0/23 Div1Int1
192.168.164.0/23 Div1Int1
192.168.166.0/24 Div1Int1
192.168.170.0/24 Div1Int1
192.168.60.0/24 Div1Int2
192.168.110.0/23 Div1Int2
192.168.120.0/23 Div1Int2
192.168.122.0/23 Div1Int2
192.168.124.0/24 Div1Int2
192.168.130.0/23 Div1Int2
192.168.140.0/23 Div1Int2
192.168.142.0/23 Div1Int2
192.168.150.0/23 Div1Int2
192.168.160.0/23 Div1Int2
192.168.162.0/23 Div1Int2
192.168.164.0/23 Div1Int2
192.168.166.0/24 Div1Int2
192.168.170.0/24 Div1Int2
192.168.121.0/24 Div1log1
192.168.122.0/24 Div1log2
192.168.123.0/24 Div1log3
192.168.124.0/24 Div1log4
192.168.141.0/24 Div1ops1
192.168.142.0/24 Div1Ext0
192.168.143.0/24 Div1Ext1
192.168.40.0/24 Div1Ext2
192.168.111.0/24 Div1Ext3
192.168.20.0/24 Div1Ext4
192.168.164.0/24 Div1Ext5
192.168.166.0/24 Div1Ext6
192.168.165.0/24 Div1Ext7
192.168.50.0/24 Div1Ext8
192.168.161.0/24 Div1Ext9
192.168.162.0/24 Div0Int0
192.168.163.0/24 Div0Int1
EOF
<2>$ rwpmapbuild --input-file=fccnets.pmap.txt \
  --output-file=fccnets.pmap
<3>$ file fccnets.pmap.txt fccnets.pmap
fccnets.pmap.txt: ASCII text
fccnets.pmap:      SiLK, PREFIXMAP v2, Little Endian, Uncompressed
```

## Predefined Prefix Maps: Country Codes and Address Types

SiLK has two predefined prefix maps to facilitate common analyses: filtering by country codes and by address types (internal, external, non-routable). They can be used as input to SiLK commands just like user-defined pmaps.

**Filtering by Country Code.** Country codes identify the nations where IP addresses are registered and are used by the Root Zone Database (e.g., see <https://www.iana.org/domains/root/db>). They are described in a `country_codes.pmap` file in the `share` directory underneath the directory where SiLK was installed or in the file specified by the `SILK_COUNTRY_CODES` environment variable.

If the current SiLK installation does not have this file, either contact the administrator that installed SiLK or look up this information on the Internet.<sup>14</sup> Country code maps are not in the normal pmap binary format and cannot be built using `rwmapbuild`.

**Filtering Internal, External, and Non-routable Addresses.** For common separation of addresses into specific types, normally internal versus external, a special pmap file may be built in the `share` directory underneath the directory where SiLK was installed. This file, `address_types.pmap`, contains a list of CIDR blocks that are labeled `internal`, `external`, or `non-routable`.

The `rwfilter` parameters `--stype` or `--dtype` use this pmap to isolate internal and external IP addresses. The `rwcut` parameter `--fields` specifies the display of this information when its argument list includes `sType` or `dType`. A value of 0 indicates `non-routable`, 1 is `internal`, and 2 is `external`. The default value is `external`.

## Using Prefix Maps to Filter Flow Records

You can use prefix maps to filter network flow records with the `rwfilter` command. This allows you to pull and partition SiLK repository data based on how it is labeled.

The pmap provides context for filtering network traffic and allows you to perform complex filtering operations in a single analytic. For example, you could create a user-defined pmap that labels IP addresses in network spaces of interest (such as the one created in Example 6.15) and filter records based on their subnetworks. You could use the predefined country code pmap (`country_codes.pmap`) to filter source or destination IP addresses based on their country of origin. You could create a user-defined port-protocol pmap to filter records for specific port and protocol combinations in order to search for unusual protocols. These types of analysis are easier to perform with data that is labeled via prefix map.

`rwfilter` supports four pmap parameters.

- `--pmap-file` specifies which compiled prefix map file to use and optionally associates a mapname with that pmap. **This switch must be specified before the other prefix map switches.**
- `--pmap-any-mapname` specifies the the set of labels used to filter records based on both source and destination IP addresses. Any source or destination IP addresses that matches an IP address with the specified label passes the filter.

---

<sup>14</sup>One such source is the GeoIP® Country database or free GeoLite™ database created by MaxMind® and available at <https://www.maxmind.com/>; the SiLK tool `rwgeoip2ccmap` converts this file to a country-code pmap.

- `--pmap-src-mapname` and `--pmap-dst-mapname` specify the set of labels for filtering by source or destination IP address, respectively. *mapname* is the name given to the pmap during construction or in `--pmap-file`. The `--pmap-file` parameter must come before any use of the *mapname* in other parameters.

Example 6.16 shows how to use the `fccnets.pmap` file from Example 6.15 to select flow records associated with web traffic from hosts on the subnetwork `Div1Int1` in the FCC network.

1. The initial call to `rwfilter` pulls all flow records from the date of interest (`--start-date=2015/06/17`) that contain inbound and outbound web traffic (`--type=inweb,outweb`).
2. The output goes to a second `rwfilter` command that uses the `--pmap-file` parameter to load the file `fccnets.pmap` and create the mapname `fccnets`. The `--pmap-any-fccnets=Div1Int1` parameter filters for all flow records whose source IP address or destination IP address matches any of the IP addresses with the label `Div1Int1` in the pmap file. These source and destination IP addresses represent traffic that flows into, within, and out of the hosts on the `Div1Int1` subnet. (Note that the `--pmap-file` parameter comes before the `--pmap-any` parameter.)
3. Records that pass the pmap-based filter are saved in `fccnets.rw`.

Example 6.16: Using Pmap Parameters with `rwfilter`

```

<1>$ rwfilter --start-date=2015/06/17 --type=inweb,outweb \
  --protocol=6 --pass=stdout \
| rwfilter stdin --pmap-file=fccnets:fccnets.pmap \
  --pmap-any-fccnets=Div1Int1 --pass=fccnets.rw
<2>$ rwfileinfo fccnets.rw --fields=count-records
fccnets.rw:
count-records          722193

```

## Displaying Prefix Values

To view the actual value of a prefix map label, use the `rwcut` command with the `--pmap-file` parameter. It takes an argument of a filename or a map name coupled to a filename with a colon. The map name typically comes from the argument for `--pmap-file`; if none is specified there, the name in the source file applies.

The `--pmap-file` parameter adds `src-mapname` and `dst-mapname` as arguments to `--fields`. Essentially, it tells `rwcut` to treat the pmap value as just another flow record field. The `--pmap-file` parameter must precede the `--fields` parameter. The two pmap fields display labels associated with the source and destination IP addresses.

Example 6.17 shows how to display the prefix labels for IP addresses in the flow record file `fccnets.rw` (created in Example 6.16). The names of the pmap and its file (`--pmap-file=fccnets:fccnets.pmap`) are specified.

## Sorting, Counting, and Grouping Records with Prefix Maps

You can also count, sort, and group flow records by prefix value. This lets you work with network data according to how it is labeled in the prefix map, which can be easier and more intuitive than some of the

Example 6.17: Viewing Prefix Map Labels with `rwcut`

```

<1>$ rwcut fccnets.rw --pmap-file=fccnets:fccnets.pmap \
  --fields=src-fccnets,sPort,dIP,dPort --num-recs=5
src-fccnets|sPort|          dIP|dPort|
Div1Int1|50373|      10.0.20.59|  80|
Div1Int1|63440|      10.0.20.59|  80|
Div1Int1|57862|      10.0.20.59|  80|
Div1Int1|62669|      10.0.20.59|  80|
Div1Int1|54211|      10.0.20.59|  80|

```

other methods for summarizing data. The `rwsort`, `rwgroup`, `rwstats`, and `rwuniq` tools all work with prefix maps. The prefix map parameters are the same as those used in the `rwcut` command and `sort`, `group`, and `count` records according to the values in the prefix map file.

Example 6.18 sorts flow records by the prefix value defined in `fccnets.pmap` for source IP addresses and bytes (`--fields=src-fccnets,bytes`). It then uses the `rwcut` command to display the pmap labels and port numbers. The first five records in the data file are displayed.

Notice that the results in Example 6.18 list `None` as their label. This is the default label from Example 6.15 that was assigned to IP addresses that are not included in the pmap. It indicates that these source IP addresses do not belong to the FCC network space. However, they exchanged traffic with hosts that are labeled as belonging to `Div1Int1` on the FCC network. Because the `rwfilter` command in Example 6.16 filtered for *all* records that contained IP addresses labeled as `Div1Int1` in the pmap, it included records where either the source or destination IP address was not part of this subnetwork.

Example 6.18: Sorting by Prefix Map Labels

```

<1>$ rwsort fccnets.rw --pmap-file=fccnets:fccnets.pmap \
  --fields=src-fccnets,bytes \
| rwcut --pmap-file=fccnets.pmap --fields=src-fccnets,sport \
  --num-recs=5
src-fccnets|sPort|
None|55862|
None|55862|
None|54392|
None|54393|
None|54394|

```

Example 6.19 shows how to count the number of records in the file `fccnets.rw` with labels defined in `fccnets.pmap`. Records without a label are listed as `None`. It also displays the destination port and the number of distinct destination IP addresses.

Again, the `rwuniq` command counts records with source or destination IP addresses that are not part of the `Div1Int1` subnet. Because these hosts communicated with hosts on the `Div1Int1` subnet, they are included in the count. The records in this case all are filtered to be TCP traffic, which makes the port numbers meaningful. TCP and UDP are the main protocols that use ports, and so are the most common ones for port/protocol pmaps. SiLK does encode ICMP type and Code into the destination port (sometimes the source port), so ICMP port/protocol pmaps are also used by some analysts.

Example 6.19: Counting Records by Prefix Map Labels

```
<1>$ rwuniq fccnets.rw --pmap-file=fccnets:fccnets.pmap \
  --fields=src-fccnets,dPort --values=flows,dIP-Distinct \
| head -n 5
src-fccnets|dPort|    Records|dIP-Distin|
  Div0Ext|52963|      8|          1|
    None|57113|      4|          1|
    None|58924|     12|          1|
  Div1Int1|55777|      2|          1|
```

### Querying Prefix Map Labels

When using prefix maps, you may need to look up which labels correspond to specific IP addresses or protocol-port pairs. Use the `rwmaplookup` command to query prefix map files—either user-defined pmaps or one of the two predefined pmaps that often are created as part of the SiLK installation (country codes and address types).

You can query a pmap with `rwmaplookup` by doing one of the following:

- specify the addresses or protocol-port pairs in a text file (the default),
- use the `--ipset-files` parameter to query the addresses in one or more IP sets,
- use the `--no-files` parameter to list the addresses or protocol-port pairs to be queried directly on the command line.

In any of these cases, one and only one of `--country-codes`, `--address-types`, or `--map-name` is used.

IP addresses are specified as described earlier in this section. For protocol-port pmaps, only the names of text files having lines in the format `protocolNumber/portNumber` or the `--no-files` parameter followed by strings in the same format are accepted. `protocolNumber` must be an integer in the range 0–255, and `portNumber` must be an integer in the range 0–65,535.

If the prefix map being queried is a protocol-port pmap, it makes no sense to query it with an IP set. `rwmaplookup` prints an error and exits if `--ipset-files` is given.

Example 6.20 shows how to use `rwmaplookup`.

- Command 1 creates a list of IP addresses and stores it in the file `ips_to_find`.
- Command 2 uses `rwmaplookup` to find the country codes associated with these addresses. If an IP address is not listed in the country code pmap, the command returns `--` as its value.
- Command 3 looks up the address types of the IP addresses. The first address has a value of 1, indicating an `internal` address. The second address has a value of 2, indicating an `external` address. The third address has a value of 0, indicating a `non-routable` address.
- Commands 4 and 5 build a protocol-port prefix map.
- Command 6 looks up protocol/port pairs from the command line.



Example 6.20: Query Addresses and Protocol/Ports with `rwpmlookup`


---

```

<1>$ cat <<END_FILE >ips_to_find
192.88.209.244
128.2.10.163
127.0.0.1
END_FILE
<2>$ rwpmlookup --country-codes ips_to_find
      key|value|
192.88.209.244|  us|
128.2.10.163|  us|
127.0.0.1|   --|
<3>$ rwpmlookup --address-types ips_to_find
      key|value|
192.88.209.244|   1|
128.2.10.163|   2|
127.0.0.1|   0|
<4>$ cat <<END_FILE >mini_icmp.pmap.txt
map-name miniicmp
default Unassigned
mode proto-port
1/0      1/0      Echo Response
1/768    1/768    Net Unreachable
1/769    1/769    Host Unreachable
1/2048   1/2048    Echo Request
END_FILE
<5>$ rwpmbuild --input-file=mini_icmp.pmap.txt --output-file=mini_icmp.pmap
<6>$ rwpmlookup --map-file=mini_icmp.pmap --no-files 1/769 1/1027
      key|          value|
1/769|Host Unreachable|
1/1027|          Unassigned|

```

---

For more information about the `rwpmlookup` command, see Appendix C.24 or enter the command `rwpmlookup --help`.

### 6.2.8 Translating IDS Signatures into `rwfilter` Calls

Traditional intrusion detection depends heavily on the presence of payloads and *signatures*: distinctive packet data that can be used to identify a particular intrusion tool. In general, the SiLK tool suite is intended for examining trends. However, it can also be used to identify specific intrusion tools. While directed intrusions are still a threat, tool-based, broad-scale intrusions are more common. Sometimes it is necessary to translate an intrusion signature into SiLK filtering rules; this section describes some standard guidelines to accomplish this task.

To convert signatures, consider the intrusion tool behavior as captured in a signature:

- What service is it targeting? This can be converted to a port number.
- What protocol does it use? This can be converted to a protocol number.
- Does it involve several protocols? Some tools, malicious and benign, will use multiple protocols, such as TCP and ICMP.
- What about packets? Buffer overflows are a depressingly common form of attack and are a function of the packet's size as well as its contents. Identifying a specific packet size can help you to figure out which intrusion tool (or tools) may have been employed for the attack.

**Hint:** When working with packet sizes, remember that the SiLK suite includes packet headers. For example, a 376-byte UDP payload will be 404 bytes long after adding 20 bytes for the IP header and eight bytes for the UDP header.

- How large are sessions? An attack tool may use a distinctive session each time (for example, a session of 14 packets with a total size of 2,080 bytes).

The `rwidsquery` command supports direct translation of rules and alert logs from the SNORT® intrusion detection system (IDS) into `rwfilter` queries. This helps an analyst examine network behavior shortly before, during, and after an alert is generated. Look for possible event triggers as well as behaviors that show the alert to be a false positive.

## Chapter 7

# Case Studies: Advanced Exploratory Analysis

This chapter features a detailed case study of exploratory analysis, using concepts from previous chapters. The study employs the SiLK workflow, SiLK tools, UNIX commands, and networking concepts to provide a practical example of exploratory analyses with network flow data.

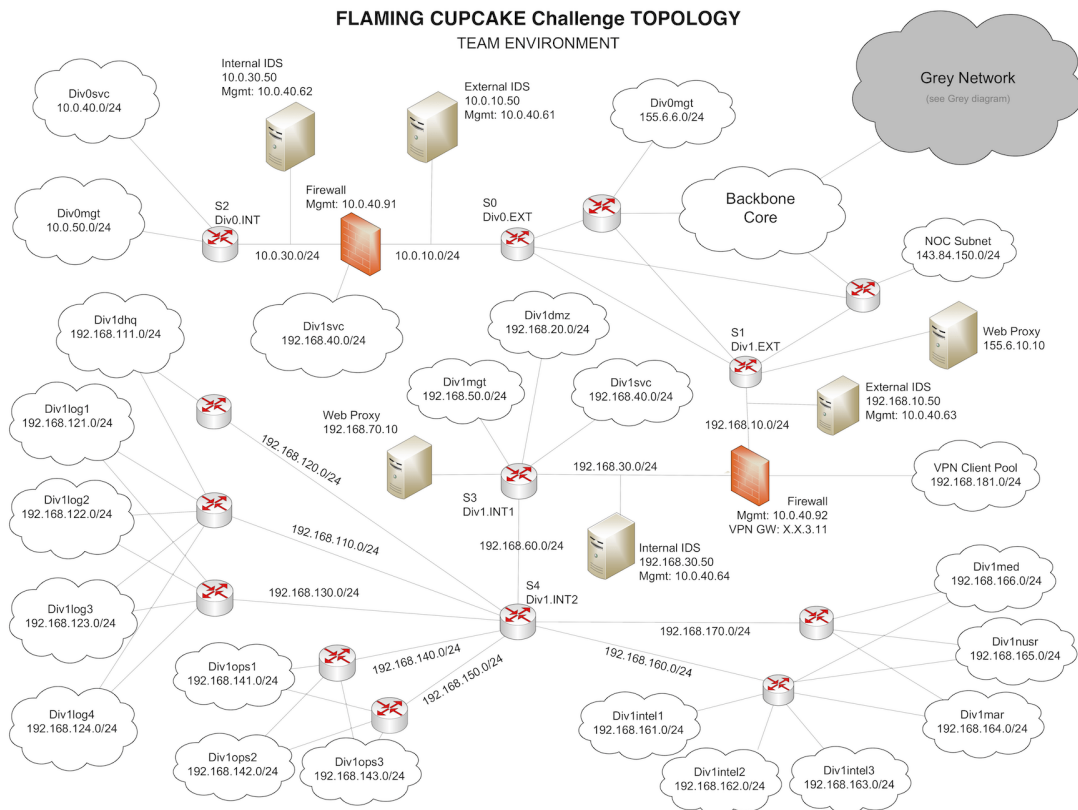
Upon completion of this chapter you will be able to

- describe how to use single-path and multi-path analyses as the building blocks of an exploratory analysis
- execute these analyses with various SiLK tools in one automated program

Each level of the exploratory analysis case is posed as a question. The case study builds upon the answers to these questions to investigate unusual network traffic and revealing changes of network behavior.

Like the previous case studies and command examples, the exploratory analysis case study uses the FCCX dataset described in Section 1.7. From the diagram in Figure 7.1, we know that sensors S0 through S4 monitor the operating network. These sensors are part of the inventory generated via the example in Chapter 5.

Figure 7.1: FCC Network Diagram



## 7.1 Level 0: Which TCP Requests are Suspicious?

In the initial phase of our exploratory analysis, we want to identify which TCP requests might represent illegitimate traffic. For most services, the flows containing client requests tend to be close in number with those containing responses to those requests. We will exploit this tendency to identify irregular traffic.

In particular, we want to look at traffic on service ports to find out which ports carry a much higher volume of inbound data than outbound data. This is one of the fingerprints (or indicators) of network scanning and several other behaviors that may be of concern. However, although this traffic imbalance is outside the range of typical behavior, it may not represent malicious activity. We need to identify these ports to take a closer look at their TCP traffic and to assess any impact it might have.

Example 7.1 shows how to find this type of data anomaly with SiLK. We will retrieve inbound and outbound TCP requests on all network ports, then find the ports that have a much higher level of inbound requests than outbound requests.

1. Command 1 uses `rwfilter` to pull inbound TCP (protocol 6) traffic for sensors S0 through S4 sent to all reserved ports (0 - 1023), which are dedicated to network services. It pipes the results to the `rwuniq` command, which counts flows and bytes for each destination port, then sorts the results to present the ports in ascending order.
2. Similarly, command 2 uses `rwfilter` to pull outbound TCP traffic and `rwuniq` to count and sort traffic for each source port.
3. Comparing the corresponding ports in the results for commands 1 and 2, we see that most of the service ports carry similar levels of inbound and outbound TCP traffic. However, ports 21, 22, and 591 carry higher inbound than outbound TCP traffic.

Activity on these ports will be investigated further.

Example 7.1: Looking for Service Ports with Higher Inbound than Outbound TCP Traffic

```

<1>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
--end=2015/06/30 --type=in --proto=6 --dport=0-1023 \
--pass=stdout \
| rwuniq --fields=dport --values=flows,bytes --sort
dPort|    Records|          Bytes|
 21|      6184|      418452|
 22|      6180|      481760|
 53|       35|       88237|
 88|     47187|     74586782|
135|      6996|     6940590|
137|      6064|     364320|
139|     57313|     93274792|
389|     22095|    118221682|
445|     81039|    363318703|
591|    112702|    76762887|
<2>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
--end=2015/06/30 --type=out --proto=6 --sport=0-1023 \
--pass=stdout \
| rwuniq --fields=sport --values=flows,bytes --sort
sPort|    Records|          Bytes|
 21|       200|       28288|
 22|       204|       99024|
 53|       35|       18860|
 88|     47184|     76069648|
135|      6991|     4052306|
137|        80|        3200|
139|     51631|     47100178|
389|     22214|    110222710|
445|     75426|    205229733|
591|     55161|    17580192|

```

## 7.2 Level 1: How Can We Identify and React to Illegitimate Requests?

The next step in our exploratory analysis is to separate normal and abnormal TCP requests on the service ports identified in Section 7.1. Specifically, we need to identify mismatched TCP flows: flows that have either no request or no response. These flows are odd, and worth examining to see if they are malicious. The goal is to describe this behavior in a way that supports further analysis.

Example 7.2 shows how to detect mismatched flows and identify the IP addresses of their sources, contrasting those addresses with sources of matched flows.

1. Command 1 uses `rwfilter` to pull inbound TCP traffic for sensors `S0` through `S4` on the suspect destination ports 21, 22, and 591. It then uses `rwsort` to sort this traffic by source IP address, destination port, destination IP address, source port, protocol, and start time. The sorted inbound flow records are saved in the file `app-in.raw`. This sort order sets up the data for matching in command 3.
2. Similarly, command 2 uses `rwfilter` to pull outbound TCP traffic for the sensors and ports of interest. It sorts this traffic by destination IP address, source port, source IP address, protocol, and start time. The sorted outbound flow records are saved in the file `app-out.raw`. This sort order allows records to be matched efficiently with those from command 1, using `rwmatch`.
3. To find mismatched flows, command 3 uses `rwmatch` to match queries in the inbound TCP flows in `app-in.raw` to responses in the outbound TCP flows in `app-out.raw`. Mismatched flows that do not belong to sessions can indicate illegitimate activity.
  - `--relate=1,2` matches the inbound source IPs to the outbound destination IPs.
  - `--relate=2,1` matches the inbound destination IPs to the outbound source IPs.
  - `--relate=3,4` matches the inbound source ports to the outbound destination ports.
  - `--relate=4,3` matches the inbound destination ports to the outbound source ports.
  - `--relate=5,5` matches the inbound and outbound protocols.
  - `--unmatched=b` saves the unmatched inbound and outbound records instead of discarding them. These unmatched records are the ones we will want to investigate for illegitimate behavior.

The matched records are indicated by setting the (mostly unused) next-hop IP address. Rather than a real IP address, `rwmatch` uses 0 followed by a positive integer value for a request, and 255 followed by the corresponding integer for a response. For a request without a response, `rwmatch` uses 0.0.0.0 for a response without a request, `rwmatch` uses 255.0.0.0. Command 3 then calls `rwsort` to sort the matched records by start time and saves them in a temporary file, `temp-match.raw`.

4. Command 4 runs `rwfilter` on `temp-match.raw` to filter records that have a next-hop IP indicating unmatched inbound requests, then saves these records in `temp-noresp.raw`. A second call to `rwfilter` filters records that have a next-hop IP indicating unmatched outbound responses, then saves those records in `temp-noreq.raw`. The records that fail both filters represent flows with matching responses and are saved in `app-match.raw`.
5. Command 5 runs the `rwstats` command on `temp-noresp.raw` to display information about the top five source IP addresses and destination ports of flows with no response flows. Since there are only two source IP addresses in this data (one with two destination ports), only three are displayed.

6. Command 6 runs the `rwstats` command on `app-match.raw` to display information about the top five source IP addresses and source ports of flows that do have matching query and response flows. Since there are 3,848 source IP addresses for the matching records, the top five are shown. None of these are the sources for the unmatched records.

Given the distribution of source addresses between the matched and unmatched traffic, these sources are clearly worth investigating further.

#### Example 7.2: Identifying Abnormal TCP Flows and their Originating Hosts

```
<1>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
--end=2015/06/30 --type=in --proto=6 --dport=21,22,591 \
--pass=stdout \
| rwsort --fields=sip,dport,dip,sport,protocol,stime \
--output=app-in.raw
<2>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
--end=2015/06/30 --type=out --proto=6 --sport=21,22,591 \
--pass=stdout \
| rwsort --fields=dip,sport,sip,dport,protocol,stime \
--output=app-out.raw
<3>$ rwmatch --relate=1,2 --relate=2,1 --relate=3,4 --relate=4,3 \
--relate=5,5 --unmatched=b app-in.raw app-out.raw stdout \
| rwsort --fields=stime --output=temp-match.raw
<4>$ rfilter temp-match.raw --next=0.0.0.0 \
--pass=temp-noresp.raw --fail=stdout \
| rfilter stdin --next=255.0.0.0 --pass=temp-noreq.raw \
--fail=app-match.raw
<5>$ rwstats --fields=sip,dport --values=flows --count=5 \
temp-noresp.raw
INPUT: 69810 Records for 3 Bins and 69810 Total Records
OUTPUT: Top 5 Bins by Records
      sIP|dPort|  Records|  %Records|  cumul_%|
      10.0.40.21| 591|    57582| 82.483885| 82.483885|
      192.168.181.8| 22|    6180|  8.852600| 91.336485|
      192.168.181.8| 21|    6048|  8.663515|100.000000|
<6>$ rwstats --fields=sip,sport --values=flows --count=5 \
app-match.raw
INPUT: 110478 Records for 3848 Bins and 110478 Total Records
OUTPUT: Top 5 Bins by Records
      sIP|sPort|  Records|  %Records|  cumul_%|
192.168.165.216| 591|    1708|  1.546009|  1.546009|
192.168.161.124| 591|    1691|  1.530621|  3.076631|
192.168.122.195| 591|    1663|  1.505277|  4.581908|
192.168.122.141| 591|    1657|  1.499846|  6.081754|
192.168.121.57| 591|    1657|  1.499846|  7.581600|
```



## 7.3 Level 2: What are the Illegitimate Sources and Destinations Doing?

In the next part of our exploratory analysis, we will investigate the activities of the illegitimate source and destination hosts identified in Section 7.2 to see what patterns emerge.

### 7.3.1 Level 2A: What are the Illegitimate Source IPs Doing?

First, we will take a look at activity on the illegitimate destination IP addresses. As shown in Example 7.3, we will find the sources of illegitimate traffic and look at the network behavior associated with scanning. Scan queries typically have low byte counts and often lack corresponding responses, since the scanned hosts lack the service being sought.

1. Command 1 calls `rwbag` and `rwbagtool --coverset` to create a set of source IPs that did not have matching queries (`noresp.set`). As input, it uses the file of no-response flows created in Section 7.2 (`temp-noresp.raw`).
2. Command 2 calls `rwfilter` to pull records coming from the source IP addresses `noresp.set`—in other words, records with source IPs that produced unmatched flows. It saves these records in `sources.raw`.
3. To find the actual scanning flows, command 3 uses `rwfilter` on `sources.raw` to filter flows with very low byte counts (0-60 bytes), indicating those with only a header, or with a header and optional extensions. It then uses the `rwuniq` command to profile these flows and saves the results in `source.txt`. This file is human readable, and contains a breakdown by protocol and bytes per flow, showing how many hosts were the destination of these flows, when the earliest of them started, and when the latest ended.

A copy of the low-volume flows is sent to a second call to `rwuniq` to save the earliest start and latest end times of the flows from each source IP to the file `source-fields.txt`. The output here is generated without column headings and vertical bar delimiter to facilitate processing by commands 5-9. As it happens, all of these low-volume flows come from a single source IP address, so the `source-fields.txt` only contains one line.

4. Command 4 displays the contents of the human-readable profile. There are several features in this output:
  - There are few distinct byte sizes in these results: only one for TCP flows, and only two for UDP flows.
  - The TCP flows are by far the most numerous, and go to by far the most distinct addresses.
  - The earliest start times are all within a one minute range.
  - The latest end times are all within a five second span.

These results support an interpretation of this behavior as scan traffic, and of a common direction behind this traffic. That all of these flows come from a single source IP address further supports this interpretation.

Example 7.3: Finding Activity of Illegitimate Destination IP Addresses

```

<1>$ rwbag --bag-file=sipv4,flows,stdout temp-noresp.raw \
| rwbagtool --coverset --output=noresp.set
<2>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
--end=2015/06/30 --type=in,out --sipset=noresp.set \
--pass=sources.raw
<3>$ rfilter sources.raw --bytes=0-60 --pass=stdout \
| rwuniq --fields=protocol,bytes \
--values=flows,distinct:dip,stime-earliest,etime-latest \
--sort --output=sources.txt --copy=stdout \
| rwuniq --fields=sip --values=stime-earliest,etime-latest \
--no-titles --delim=' ' --output=source-fields.txt
<4>$ cat sources.txt
pro|      bytes|  Records|dIP-Distin|      sTime-Earliest|      eTime-Latest|
 6|         60|   29492|    768|2015/06/17T16:12:58|2015/06/17T16:41:21|
17|         29|     68|     17|2015/06/17T16:13:54|2015/06/17T16:41:21|
17|         42|    136|     17|2015/06/17T16:13:54|2015/06/17T16:41:26|
<5>$ srcArray=( $(cat source-fields.txt) )

```

### 7.3.2 Level 2B: What Behavior Changes do Destination IPs Show?

Next, we will investigate traffic patterns on the destination hosts (the targets of the scans). Using the start times and end times of the scan from Example 7.3, we will look at traffic patterns before and afterwards as shown in Example 7.4.

1. Commands 1 through 5 locate the start times (`StTime`, `StEpoch`) and end times (`EnTime`, `EnEpoch`), using the `source-fields.txt` file created in Example 7.3. This identifies the “before” and “after” boundaries of the scan.
  - Command 1 stores the contents of the file (one line) in a shell array for ease of reference to the fields.
  - Command 2 converts the earliest start time of the scan into an integer epoch value (number of seconds since Midnight January 1, 1970). This format is used to make it easy to calculate.
  - Command 3 subtracts one from the epoch value to get a time briefly before the scan activity started, then uses the string processing language `awk` to convert the epoch date back into a SiLK formatted date.
  - Commands 4 and 5 do an analogous process to commands 2 and 3, but with the ending time of the scanning activity, producing a value just after the scanning ended.
2. Command 6 uses `rfilter` to pull inbound and outbound non-web traffic for the destination IPs in `noresp.set` in the time window *before* the start of the scan, saving these flows to `dest-before.raw`.
3. Command 7 uses `rfilter` to pull inbound and outbound traffic for the illegitimate IPs in `noresp.set` in the time window *after* the end of the scan, saving these files to `dest-after.raw`.
4. For clarity of display, command 8 calls `rwuniq`, then uses the `head` command to pull off the column headers and store them in the file `myhead`.

Example 7.4: Finding Changed Behavior in Destination IPs

```

<1>$ srcArray=( $(cat source-fields.txt) )
<2>$ StEpoch=$( date -d ${srcArray[1]} +"%s")
<3>$ StTime=$(echo $(( $StEpoch - 1 )) | awk '{print \
    strftime("%Y/%m/%dT%T", $1)}')
<4>$ EnEpoch=$( date -d ${srcArray[2]} +"%s")
<5>$ EnTime=$(echo $(( $EnEpoch + 1 )) | awk '{print \
    strftime("%Y/%m/%dT%T", $1)}')
<6>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
    --end=2015/06/30 --type=in,out --dipset=noresp.set \
    --etime=2015/06/01-2015/06/17T09:12:57 \
    --pass=dest-before.raw
<7>$ rfilter --sensor=S0,S1,S2,S3,S4 --start=2015/06/01 \
    --end=2015/06/30 --type=in,out --dipset=noresp.set \
    --stime=2015/06/17T09:41:27-2015/06/30 \
    --pass=dest-after.raw
<8>$ rwuniq --fields=bytes,protocol \
    --values=Flows,distinct:sip,distinct:dip dest-before.raw \
    | head -1 >myhead
<9>$ cat myhead; \
rwuniq --fields=bytes,protocol \
    --values=Flows,distinct:sip,distinct:dip --sort \
    dest-before.raw \
| tail -5
    bytes|pro|    Records|sIP-Distin|dIP-Distin|
    884| 17|        2|        1|        1|
    988|  1|        4|        1|        1|
   1028| 17|        5|        2|        1|
   1326| 17|        1|        1|        1|
   1976|  1|        3|        1|        1|
<10>$ cat myhead; \
rwuniq --fields=bytes,protocol \
    --values=Flows,distinct:sip,distinct:dip --sort \
    dest-after.raw \
| tail -5
    bytes|pro|    Records|sIP-Distin|dIP-Distin|
   5128|  1|        8|        2|        1|
   5488|  1|        4|        1|        1|
   5756|  1|        4|        1|        1|
   5844|  1|        4|        1|        1|
   6020|  1|        4|        1|        1|

```

5. After displaying the column headers with `cat`, command 9 uses `rwuniq` to calculate the counts of flows for each byte size in `dest-before.rw`, which contains data from the time period before the start of scanning. It sorts the byte sizes in ascending order, then calls `tail` to display the largest flows. (Because we use `tail` for this, we had to separately save and display the column headers; without that, no column labels would be shown.)
6. To profile network behavior after scanning, command 10 uses `rwuniq` to find the flows with the largest byte sizes in the file `dest-after.rw`, which contains data from the time period after scanning ended.

## 7.4 Level 3: What are the Commonalities Across The Cases?

Our exploratory analysis of the network traffic shows likely scanning behavior during a tight time frame: a 28-minute interval, ending within a few seconds.

A further look at the results reveals a definite change in behavior of the suspected scanners. One IP address is active before and during the scan. Another is active after, with much larger flows to the destination after the scan completed. All of this is highly suspicious. Even more concerning, the large traffic after the scan is all ICMP traffic (which is normally quite modest in size). These large ICMP flows are highly unusual for any benign purpose.

Future areas for investigation include

- looking for additional hosts that exhibit behavior that is similar to the scan/exploit hosts, during the periods before and after the scan
- investigating the behavior of the scan and exploit hosts for further confirmation of their malicious character
- looking for other activities associated with the scan/exploit hosts. What else are they up to?

As the analysis continues, these areas will likely suggest others to be explored. One difficulty here is knowing when to stop. Analysts need to keep the desired level of output firmly in mind, steering their explorations to provide suitable results. They need to stop either when those results are found (a likely compromise is identified, a sufficient understanding of the service has resulted, or interactions across the network are understood), or when it is clear no such results will emerge (i.e., everything is benign).

In this case, having identified a definite change in behavior, the analysis has reached its end. Information about the affected hosts could then be passed to incident handlers or system administrators for response.

## Chapter 8

# Extending the Reach of SiLK with PySiLK

This chapter discusses how to use PySiLK, the SiLK Python extension, to support analyses that are difficult to implement within the normal constraints of the SiLK tool suite. Sometimes, an analyst needs to use parts of the SiLK suite’s native functionality in a modified way. The capabilities of PySiLK simplify these analyses.

This chapter does not discuss the issues involved in composing new PySiLK scripts or how to code in the Python programming language. Several example scripts are shown, but the detailed design of each script will not be presented here.

Upon completion of this chapter, you will be able to

- explain the purpose of PySiLK in analysis
- use and modify PySiLK plug-ins to match records in `rwfilter`
- use and modify PySiLK plug-ins to add fields for `rwcut` and `rwsort`
- use and modify PySiLK plug-ins to add key fields and summary values for `rwuniq` and `rwstats`

Additional PySiLK and Python programming language resources include

- a brief guide to coding PySiLK plug-ins, provided by the `silkpython` manual page (see `man silkpython` or Section 3 of *The SiLK Reference Guide* at <https://tools.netsa.cert.org/silk/reference-guide.pdf>)
- detailed descriptions of the PySiLK structures, provided in *PySiLK: SiLK in Python* (<https://tools.netsa.cert.org/silk/pysilk.pdf>)
- larger PySiLK examples, provided in the PySiLK “tooltips” webpage (<https://tools.netsa.cert.org/confluence/display/tt/Writing+PySiLK+scripts>)
- generic programming in the Python programming language, as described in many locations on the web, particularly on the Python official website (<https://www.python.org/doc/>)

## 8.1 Using PySiLK

PySiLK is an extension to the SiLK tool suite that expands its functionality via scripts written in the Python programming language. The purpose of PySiLK is to support analytical use cases that are difficult to express, implement, and support with the capabilities natively built into SiLK, while using those capabilities where appropriate.

To extend the SiLK tools with PySiLK, first write a Python file that calls Python functions defined in the `silk.plugin` Python module. To use this Python file, specify the `--python-file` switch for one of the SiLK tools that supports PySiLK. (The `rwfilter`, `rwstats`, `rwuniq`, `rwcut`, and `rwsort` tools can all make use of PySiLK extensions.) The tool then loads the Python file and makes the new functionality available.

### 8.1.1 PySiLK Requirements

To use PySiLK:

1. Install the appropriate version of the Python language<sup>15</sup> on your system.
2. Load the PySiLK library (a directory named `silk`) in the `site-packages` directory of the Python installation.
3. To ensure that the PySiLK library works properly, set the `PYTHONPATH` environment variable to include the `site-packages` directory.

### 8.1.2 PySiLK Scripts and Plug-ins

PySiLK code comes in two forms: standalone Python programs and plug-ins for SiLK tools. Both of these forms use a Python module named `silk` that is provided by CERT<sup>®</sup> as part of the PySiLK library. PySiLK provides the capability to manipulate SiLK objects (flow records, IPsets, bags, etc.) with Python code.

For analyses that will not be repeated often or that are expected to be modified frequently, the relative brevity of PySiLK renders it an efficient alternative. As with all programming, analysts need to use algorithms that meet the speed and space constraints of the project. Often (but not always), building upon the processing of the SiLK tools by use of a plug-in yields a more suitable solution than developing a stand-alone script.

PySiLK plug-ins for SiLK tools use an additional component, also provided by CERT<sup>®</sup>, called `silkpython`. The `silkpython` component creates the application programming interfaces (APIs), simple and advanced, that connect a plug-in to SiLK tools. Currently, `silkpython` supports the following SiLK tools: `rwfilter`, `rwcut`, `rwgroup`, `rwsort`, `rwstats`, and `rwuniq`.

- For `rwfilter`, `silkpython` permits a plug-in to provide new types of partitioning criteria.
- For `rwcut`, plug-ins can create new flow record fields for display.
- For `rwgroup` and `rwsort`, plug-ins can create fields to be used as all or part of the key for grouping or sorting records.
- For `rwstats` and `rwuniq`, two types of fields can be defined: *key* fields used to categorize flow records into bins and *summary value* fields used to compute a single value for each bin from the records in those bins.

---

<sup>15</sup>Python 2.7.x for SiLK Version 3.8.

- For all of the tools, `silkpython` allows a plug-in to create new SiLK tool switches (parameters) to modify the behavior of the aforementioned partitioning criteria and fields.

The `silkpython` module provides only one function for establishing (registering) partitioning criteria (filters) for `rwfilter`. The simple API provides four functions for creating key fields for the other five supported SiLK tools and three functions for creating summary value fields for `rwstats` and `rwuniq`. The advanced API provides one function that can create either key fields or summary value fields, and permits a higher degree of control over these fields.

## 8.2 Extending `rwfilter` with PySiLK

PySiLK extends the capabilities of `rwfilter` by letting the analyst create new methods for partitioning flow records.<sup>16</sup>

For a single execution of `rwfilter`, PySiLK is much slower than using a combination of `rwfilter` parameters and usually slower than using a C-language plug-in. However, there are several ways in which using PySiLK can replace a series of several `rwfilter` executions with a single execution and ultimately speed up the overall process.

Without PySiLK, `rwfilter` has limitations using its built-in partitioning parameters:

- Each flow record is examined without regard to other flow records. That is, no state is retained.
- There is a fixed rule for combining partitioning parameters: Any of the alternative values within a parameter satisfies that criterion (i.e., the alternatives are joined implicitly with a logical *or* operation). All parameters must be satisfied for the record to pass (i.e., the parameters are joined implicitly with a logical *and* operation).
- The types of external data that can assist in partitioning records are limited. IP sets, tuple files, and prefix maps are the only types provided by built-in partitioning parameters.

PySiLK is useful to expand the capabilities of `rwfilter` in these cases:

- Information from prior records may help to partition subsequent records into the pass or fail categories.
- A series of nontrivial alternatives form the partitioning condition.
- The partitioning condition employs a control structure or data structure.

### 8.2.1 Using PySiLK to Incorporate State from Previous Records: Eliminating Inconsistent Sources

For an example of where some information (or *state*) from prior records may help in partitioning subsequent records, consider Example 8.1. This script (`ThreeOrMore.py`) passes all records that have a source IP address used in two or more prior records. This can be useful if you want to eliminate casual or inconsistent sources of particular behavior. The `addrRefs` variable is the record of how many times each source IP address has

Example 8.1: ThreeOrMore.py: Using PySiLK for Memory in `rwfilter` Partitioning

---

```

import sys    # stderr

bound = 3    # default threshold for passing record
addrRefs={} # key = IP address, value = reference count

def threeOrMore(rec):
    global addrRefs # allow modification of addrRefs

    keyval = rec.sip # change this to count on different field
    addrRefs[keyval] = addrRefs.get(keyval, 0) + 1
    return addrRefs[keyval] >= bound

def set_bound(integer_string):
    global bound

    try:
        bound = int(integer_string)
    except ValueError:
        print >>sys.stderr, '--limit value, %s, is not an integer.' % integer_string

def output_stats():
    AddrsWithEnufFlows = len([1 for k in addrRefs.keys()
                              if addrRefs[k] >= bound])
    print >>sys.stderr, 'SIPs: %d; SIPs meeting threshold: %d' % (len(addrRefs),
                                                                    AddrsWithEnufFlows)

register_filter(threeOrMore, finalize=output_stats)
register_switch('limit', handler=set_bound, help='Threshold for passing')

```

---



been seen in prior records. The `threeOrMore` function holds the Python code to partition the records. If it determines the record should be passed, it returns `True`; otherwise it returns `False`.

In Example 8.1, the call to `register_filter` informs `rwfilter` (through `silkpython`) to invoke the specified Python function (`threeOrMore`) for each flow record that has passed all the built-in `SiLK` partitioning criteria. In the `threeOrMore` function, the `addrRefs` dictionary is a container that holds entries indexed by an IP address and whose values are integers.

When the `get` method is applied to the dictionary, it obtains the value for the entry with the specified key, `keyval`, if such an entry already exists. If this is the first time that a particular key value arises, the `get` method returns the supplied default value, zero. Either way, one is added to the value obtained by `get`. The `return` statement compares this incremented value to the `bound` threshold value and returns the Boolean result to `silkpython`, which informs `rwfilter` whether the current flow record passes the partitioning criterion in the `PySiLK` plug-in.

In Example 8.1, the `set_bound` function is not required for the partitioning to operate. It provides the capability to modify the threshold that the `threeOrMore` function uses to determine which flow records pass. The call to `register_switch` informs `rwfilter` (through `silkpython`) that the `--limit` parameter is acceptable in the `rwfilter` command after the `--python-file=ThreeOrMore.py` parameter, and that if the user specifies the parameter (e.g., `--limit=5`) the `set_bound` function will be run to modify the `bound` variable before any flow records are processed. The value provided in the `--limit` parameter will be passed to the `set_bound` function as a string that needs to be converted to an integer so it can participate later in numerical comparisons.

If the user specifies a string that is not a representation of an integer, the conversion will fail inside the `try` statement, raising a `ValueError` exception and displaying an error message; in this case, `bound` is not modified.

### 8.2.2 Using `PySiLK` to Incorporate State from Previous Records: Detecting Port Knocking

For an example in which some information (or *state*) from prior records may help in partitioning subsequent records, consider *port knocking*. This is a technique used to thwart port scanning. Port scanning involves sending single packets to particular ports on a target host to see what response, if any, is returned by the target. The response, or lack of one, is interpreted to determine if there is a service available on that port on the target host. Port knocking is employed by the administrator of the target host to make all ports look as if there are no services available on any of them.

Port knocking requires legitimate users (or their software) to know the secret combination of actions that must be taken before an attempt is made to connect to a service port. These actions consist of attempts to connect (or *knocks*) to certain other ports in the correct order right before attempting a connection to the service port. Achieving the correct sequence of knocks creates a temporary rule in the firewall to allow the sender of the knocks to connect to a particular service. Profiling a network to obtain situational awareness could include port knocking detection to explain what would otherwise look like strange traffic.

Example 8.2 implements a plug-in for `rwfilter` that takes a simple approach. The plug-in requires that its input be sorted by IP addresses and time. That way the port knocks appear in the input right before the service connection attempt. This means that the plug-in needs to retain state for only one connection attempt at a time, simplifying the code and greatly reducing the memory requirements.

---

<sup>16</sup>These partitioning methods may also calculate values not normally part of `rwfilter`'s output, typically emitting those values to a file or to the standard error stream to avoid conflicts with flow record output.

The plug-in looks for three consecutive flow records where the first two (the port knocks) attempt to initiate TCP connections to different ports, but there are no following packets with the ACK flag that would indicate the connection had been established. After the two port knock flows, there must be a flow for a third port which does have additional packets with the ACK flag.

These criteria are somewhat simple. We could add constraints such as specifying that the three ports must have a certain ordinal relationship (e.g., lower-higher-lower). However, Example 8.2 shows the essential elements. When the three-flow sequence is found, the third flow passes the filter and a text record is displayed.

Example 8.2: portknock.py: Using PySiLK to Retain State in `rwfilter` Partitioning

```
import datetime # timedelta()
import sys # stdout

REQDKNOCKS = 2 # number of required knocks with distinct port numbers
INTERVAL = datetime.timedelta(seconds=5) # knocks & conn this close in time
TCP = 6 # protocol number

portListWidth = REQDKNOCKS * 7
lastsip = None

def note_first_knock(rec):
    global lastsip, lastdip, portlist, lastetime
    lastsip = rec.sip
    lastdip = rec.dip
    portlist = [rec.dport]
    lastetime = rec.etime
    return

def examine_flow(rec):
    global lastsip, lastdip, portlist, lastetime
    if (rec.protocol == TCP and rec.initial_tcpflags is not None and
        rec.initial_tcpflags.matches('S/SA')): # initial SYN (client to server)
        if lastsip is not None and rec.sip == lastsip and rec.dip == lastdip:
            if rec.stime - lastetime <= INTERVAL:
                if rec.session_tcpflags.ack: # established connection
                    # connected to knocked port or insufficient knocks?
                    if rec.dport in portlist or len(portlist) < REQDKNOCKS:
                        lastsip = None
                else: # enough prior knocks?
                    sys.stdout.write('%15s %15s %*s %5d\n' % (lastsip, lastdip,
                                                                portListWidth, portlist, rec.dport))
                    lastsip = None
                    return True # flow record passes filter
            else: # connection not established; just a knock
                if rec.dport in portlist: # already seen this port
                    note_first_knock(rec) # start over as 1st knock
                else:
                    if len(portlist) >= REQDKNOCKS: # add a knock
                        del portlist[0] # delete oldest knock, make room for new
                        portlist.append(rec.dport)
                        lastetime = rec.etime
        # last knock was too long ago
```

```

        elif rec.session_tcpflags.ack: # established connection
            lastsip = None
        else: # too long ago and connection not established; just a knock
            note_first_knock(rec) # start over as 1st knock
        # new sip and/or dip
        elif not rec.session_tcpflags.ack: # conn not established; just a knock
            note_first_knock(rec) # start over as 1st knock
    return False # flow record fails filter

def show_heading():
    sys.stdout.write('%15s %15s %*s %5s\n' % ('sIP', 'dIP', portListWidth,
                                             'Knock-Ports', 'Estab'))

register_filter(examine_flow, initialize=show_heading)

```

### 8.2.3 Using PySiLK with rfilter in a Distributed or Multiprocessing Environment

An analyst could use a PySiLK script with `rfilter` by first calling `rfilter` to retrieve the records that satisfy a given set of conditions, then piping those records to a second `rfilter` call that uses the `--python-file` parameter to invoke the script. This is shown in Example 8.3. This syntax is preferred to simply including the `--python-file` parameter on the first call, since its behavior is more consistent across execution environments. If `rfilter` is running on a multiprocessor configuration, running the script on the first `rfilter` call cannot be guaranteed to behave consistently for a variety of reasons, so running PySiLK scripts via a piped `rfilter` call is more consistent.

Example 8.3: Calling `ThreeOrMore.py`

```

<1>$ rfilter --start-date=2015/06/02 --end-date=2015/06/18 \
  --type=inweb --protocol=6 --dport=443 \
  --bytes-per-packet=65- --packets=4- \
  --flags-all=SAF/SAF,SAR/SAR --pass=stdout \
| rfilter stdin --python-file=ThreeOrMore.py \
  --pass=web.rw
SIPs: 81; SIPs meeting threshold: 81

```

### 8.2.4 Simple PySiLK with rfilter --python-expr

Some analyses that do not lend themselves to solutions with just the SiLK built-in partitioning parameters may be so simple with PySiLK that they center on an expression that evaluates to a Boolean value. Using the `rfilter --python-expr` parameter will cause `silkpython` to provide the rest of the Python plug-in program.

Example 8.4 partitions flow records that have the same port number for their source port and destination port (`sport` and `dport`). Although the name for the flow record object is specified by a function parameter in user-written Python files, with `--python-expr`, the record object is always called `rec`. The source port

---

 Example 8.4: Using `--python-expr` for Partitioning
 

---

```
<1>$ rfilter flows.rw --protocol=6,17 --python-expr='rec.sport==rec.dport' \
      --pass=equalports.rw
```

---

and destination port therefore can be specified as `rec.sport` and `rec.dport`. Checking whether their values are equal becomes very simple.

With `--python-expr`, it is not possible to retain state from previous flow records as in Example 8.1. Nor is it possible to incorporate information from sources other than the flow records. Both of these require a plug-in invoked by `--python-file`.

## 8.2.5 PySiLK with Complex Combinations of Rules

Example 8.5 shows an example of using PySiLK to filter for a condition with several alternatives. This code is designed to identify virtual private network (VPN) traffic in the data, using IPsec, OpenVPN<sup>®</sup>, or VPNz<sup>®</sup>. This involves having several alternatives, each matching traffic either for a particular protocol (50 or 51) or for particular combinations of a protocol (17) and ports (500, 1194, or 1224). This could be done using a pair of `rfilter` calls (one for UDP [17] and one for both ESP [50] and AH [51]) and `rwcat` to put them together, but this is less efficient than using PySiLK.

---

 Example 8.5: `vpn.py`: Using PySiLK with `rfilter` for Partitioning Alternatives
 

---

```
def vpnfilter(rec):
    return ((rec.protocol == 17 and # UDP
            (rec.dport in (500, 1194, 1224) or # IKE, OpenVPN, VPNz
             rec.sport in (500, 1194, 1224) ) )
           or rec.protocol in (50, 51) ) # ESP, AH

register_filter(vpnfilter)
```

---

## 8.2.6 Use of Data Structures in Partitioning

Example 8.6 shows the use of a data structure in an `rfilter` condition. This particular case identifies internal IP addresses responding to contacts by IP addresses in certain external blocks. The difficulty is that the response is unlikely to go back to the contacting address and likely instead to go to another address on the same network. Matching this with conventional `rfilter` parameters is very slow and repetitive. By building a list of internal IP addresses and the networks they've been contacted by, `rfilter` can partition records based on this list using the PySiLK script in Example 8.6, called `matchblock.py`.

In Example 8.6, lines 1 and 2 import objects from two modules. Line 3 sets a constant (with a name in all uppercase by convention). Line 4 creates a global variable to hold the name of the file containing external netblocks and gives it a default value. Lines 6, 10, and 32 define functions to be invoked later. Line 42 informs `silkpython` of two things: (1) that the `open_blockfile` function should be invoked after all command-line switches (parameters) have been processed and before any flow records are read and (2) that in addition to any other partitioning criteria, every flow record must be tested with the `match_block` function to determine

if it passes or fails. Line 43 tells `silkpython` that `rwfilter` should accept a `--blockfile` parameter on the command line and process its value with the `change_blockfile` function before the initialization function, `open_blockfile`, is invoked.

When `open_blockfile` is run, it builds a list of external netblocks for each specified internal address. Line 25 converts the specified address to a PySiLK address object; if that's not possible, a `ValueError` exception is raised, and that line in the blockfile is skipped. Line 26 similarly converts the external netblock specification to a PySiLK IP wildcard object; if that's not possible, a `ValueError` exception is raised, and that line in the file is skipped. Line 26 also appends the netblock to the internal address's list of netblocks; if that list does not exist, the `setdefault` method creates it.

When each flow record is read by `rwfilter`, `silkpython` invokes `match_block`, which tests every external netblock in the internal address's list to see if it contains the external, destination address from the flow record. If an external address is matched to a netblock in line 35, the test passes. If no netblocks in the list match, the test fails in line 39. If there is no list of netblocks for an internal address (because it was not specified in the blockfile), the test fails in line 38.

Example 8.7 uses command-line parameters to invoke the Python plug-in and pass information to the plug-in script (specifically the name of the file holding the block map). Command 1 displays the contents of the block map file. Each line has two fields separated by a comma. The first field contains an internal IP address; the second field contains a wildcard expression (which could be a CIDR block or just a single address) describing an external netblock that has communicated with the internal address. Command 2 then invokes the script using the syntax introduced previously, augmented by the new parameter.

Example 8.6: matchblock.py: Using PySiLK with rwfilter for Structured Conditions

```

from silk import IPAddr,IPWildcard
2 import sys # exit(), stderr
  PLUGIN_NAME = 'matchblock.py'
  blockname='blocks.csv'
5
def change_blockfile(block_str):
  global blockname
8   blockname = block_str

def open_blockfile():
11  global blockfile, blockdict
    try:
      blockfile = open(blockname)
14  except IOError, e_value:
      sys.exit('%s: Block file: %s' % (PLUGIN_NAME, e_value))
    blockdict = dict()
17  for line in blockfile:
      if line.lstrip()[0] == '#': # recognize comment lines
          continue # skip entry
20  fields = line.strip().split(',') # remove NL and split fields on commas
      if len(fields) < 2: # too few fields?
          print >>sys.stderr, '%s: Too few fields: %s' % (PLUGIN_NAME, line)
23  continue # skip entry
      try:
          idx = IPAddr(fields[0].rstrip())
26  blockdict.setdefault(idx, []).append(IPWildcard(fields[1].strip()))
      except ValueError: # field cannot convert to IPAddr or IPWildcard
          print >>sys.stderr, '%s: Bad address or wildcard: %s' % (PLUGIN_NAME, line)
29  continue # skip entry
    blockfile.close()

32 def match_block(rec):
    try:
      for netblock in blockdict[rec.sip]:
35  if rec.dip in netblock:
          return True
    except KeyError: # no such inside addr
38  return False
    return False # no netblocks match

41 # M A I N
register_filter(match_block, initialize=open_blockfile)
register_switch('blockfile', handler=change_blockfile,
44 help='Name of file that holds CSV block map. Def. blocks.csv')

```

Example 8.7: Calling `matchblock.py`


---

```
<1>$ cat blockfile.csv
198.51.100.17, 192.168.0.0/16
203.0.113.178, 192.168.x.x
<2>$ rwwfilter out_month.rw --protocol=6 --dport=25 --pass=stdout \
  | rwwfilter stdin --python-file=matchblock.py \
    --blockfile=blockfile.csv --print-statistics
Files      1.  Read      375567.  Pass          8.  Fail      375559.
```

---

### 8.3 Extending SiLK with Fields Defined with PySiLK

Five SiLK tools support fields defined with PySiLK: `rwwcut`, `rwwgroup`, `rwwsort`, `rwwstats`, and `rwwuniq`. All five support newly defined *key* fields, although for `rwwcut` these fields are not really the key to any sorting or grouping operation; they are simply available for display. Two of the tools, `rwwstats` and `rwwuniq`, support newly defined *summary value* fields. For fields that require the use of the advanced API, the programmer will want to determine which of the five tools will be used with these fields to avoid unnecessary programming. By examining the characteristics of the new field the programmer can determine which tools can take advantage of the field.

`rwwcut` can make use of any key field that produces character strings (text). It does not matter if field values can be used in computations, comparisons, or sorting. As long as the field can be represented as text, `rwwcut` can use it. When registering such fields, a function to produce a text value must be provided.

For `rwwgroup` and `rwwsort` to utilize a field defined with PySiLK, the field registration must provide a function that produces binary (non-text) values that can be compared. For `rwwgroup`, the comparison is only for equality to determine if two consecutive records belong in the same group. For `rwwsort` the comparison must also determine the order of unequal values.

For `rwwstats` and `rwwuniq`, a key field not only needs a binary representation for grouping purposes, it also needs a function to convert the binary key (bin) to text for display purposes. Furthermore, it must make sense semantically for the field to produce a many-to-one mapping from its inputs to its value. This is necessary for the field to make a good key (or partial key) for bins. A field makes a poor binning key if nearly every record produces a unique field value, or if nearly every record produces the same field value.

### 8.4 Extending `rwwcut` and `rwwsort` with PySiLK

PySiLK is useful with `rwwcut` and `rwwsort` in these cases:

- An analysis requires a value based on a combination of fields, possibly from a number of records.
- An analyst chooses to use a function on one or more fields, possibly conditioned by the value of one or more fields. The function may incorporate data external to the records (e.g., a table of header lengths).

### 8.4.1 Computing Values from Multiple Records

Example 8.8 shows the use of PySiLK to calculate a value from the same field of two different records in order to provide a new column to display with `rwcut`. This particular case, which will be referred to as `delta.py`, introduces a `delta_msec` column, with the difference between the start time of two successive records. Potential uses for this column including ready identification of flows that occur at very stable intervals, such as keep-alive traffic or beaconing.

The plug-in uses global variables to save the IP addresses and start time between records and then returns to `rwcut` the number of milliseconds between start times. The `register_int_field` call allows the use of `delta_msec` as a new field name and gives `rwcut` the information that it needs to process the new field.

Example 8.8: `delta.py`

---

```

last_sip = None

def compute_delta(rec):
    global last_sip, last_dip, last_time
    if last_sip is None or rec.sip != last_sip or rec.dip != last_dip:
        last_sip = rec.sip
        last_dip = rec.dip
        last_time = rec.stime_epoch_secs
        deltamsec = 0
    else: # sip and dip same as previous record
        deltamsec = int(1000. * (rec.stime_epoch_secs - last_time))
        last_time = rec.stime_epoch_secs
    return deltamsec

register_int_field ('delta_msec', compute_delta, 0, 4294967295)
#                 fieldname      function      min      max

```

---

To use `delta.py`, Example 8.9 sorts the flow records by source address, destination address, and start time after pulling them from the repository. After sorting, the example passes them to `rwcut` with the `--python-file=delta.py` parameter before the `--fields` parameter so that the `delta_msec` field name is defined. Because of the way the records are sorted, if the source or destination IP addresses are different in two consecutive records, the latter record could have an earlier `sTime` than the prior record. Therefore, it makes sense to compute the time difference between two records only when their source addresses match and their destination addresses match. Otherwise, the delta should display as zero.

### 8.4.2 Computing a Value Based on Multiple Fields in a Record

Example 8.10 shows the use of a PySiLK plug-in for both `rwsort` and `rwcut` that supplies a value calculated from several fields from a single record. In this example, the new value is the number of bytes of payload conveyed by the flow. The number of bytes of header depends on the version of IP as well as the Transport-layer protocol being used (IPv4 has a 20-byte header, IPv6 has a 40-byte header, and TCP adds 20 additional bytes, while UDP adds only 8 and GRE [protocol 47] only 4, etc.).

The `header_len` variable holds a mapping from protocol number to header length. Protocols omitted from the mapping contribute zero bytes for the Transport-layer header. This is then multiplied by the number of packets and subtracted from the flow's byte total. This code assumes no packet fragmentation is occurring.



Example 8.9: Calling delta.py

```

<1>$ rfilter --type=out --start-date=2015/06/02 \
  --end-date=2015/06/18 --protocol=17 --packets=1 \
  --pass=stdout \
| rwsort --fields=sIP,dIP,sTime \
| rwcut --python-file=delta.py \
  --fields=sIP,dIP,sTime,delta_msec --num-recs=20
  sIP|          dIP|          sTime|delta_msec|
10.0.20.58| 128.8.10.90|2015/06/17T20:50:50.725| 0|
10.0.20.58| 128.8.10.90|2015/06/17T20:50:50.725| 0|
10.0.20.58| 128.8.10.90|2015/06/17T20:50:50.725| 0|
10.0.20.58| 199.7.83.42|2015/06/17T20:50:46.717| 0|
10.0.20.58| 199.7.83.42|2015/06/17T20:50:46.717| 0|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:08.030| 0|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:08.268| 237|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:08.580| 312|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:08.580| 0|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:08.674| 94|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:09.015| 341|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:09.043| 27|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:09.215| 171|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:09.397| 182|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:10.228| 830|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:10.620| 391|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:10.622| 2|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:10.622| 0|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:11.069| 447|
10.0.40.20| 10.0.20.58|2015/06/16T12:48:11.415| 345|

```

The same function is used to produce both a value for `rwsort` to compare and a value for `rwcut` to display, as indicated by the `register_int_field` call.

Example 8.10: `payload.py`: Using PySiLK for Conditional Fields with `rwsort` and `rwcut`

```

#           ICMP IGMP IPv4   TCP   UDP   IPv6  RSVP
header_len={1:8, 2:8, 4:20, 6:20, 17:8, 41:40, 46:8,
            47:4, 50:8, 51:12, 88:20, 132:12}
#           GRE   ESP    AH    EIGRP  SCTP

def bin_payload(rec):
    transport_hdr = header_len.get(rec.protocol, 0)
    if rec.is_ipv6():
        ip_hdr = 40
    else:
        ip_hdr = 20
    return rec.bytes - rec.packets * (ip_hdr + transport_hdr)

register_int_field('payload', bin_payload, 0, (1 << 32) - 1)
#           fieldname      function  min      max

```

Example 8.11 shows how to use Example 8.10 with both `rwsort` and `rwcut`. The records are sorted into payload-size order and then output, showing both the bytes and payload values.

### 8.4.3 Defining a Character String Field for `rwcut`

PySiLK has no function in the simple API for creating character-string fields. Do not be tempted to use an enumeration where there are many possible strings produced for the field; an enumeration can be quite memory-intensive. It will not sort correctly in `rwsort` or in `rwuniq` with the `--sort-output` parameter.

Creating a character-string field with the advanced API (`register_field` function) is not difficult. Starting with an example of a string field that works only with `rwcut`, Example 8.12 is a PySiLK plug-in that makes large values in the built-in `duration` field more understandable by breaking it down into days, hours, minutes, and seconds (including milliseconds). This field does not provide a binary value, so the field has no usefulness with `rwsort` or `rwgroup`, which produce non-text output. Since this field embodies the same information as the built-in `duration` field, the built-in field is better suited for use with these tools as it will yield much better performance. The usefulness of `decode_duration` with `rwstats` and `rwuniq` is dubious as well; although these tools produce textual output, the lack of a many-to-one mapping makes this field an ideal candidate for use only with `rwcut`.

Example 8.13 shows both the built-in `duration` field and the associated `decode_duration` field for several flow records.

### 8.4.4 Defining a Character String Field for Five SiLK Tools

A plug-in to create a character-string field not only for `rwcut`, but also for `rwgroup`, `rwsort`, `rwstats`, and `rwuniq` needs a little more code. In Example 8.14 a regular expression is used to provide a pattern for the site-name portion of a sensor name. Since the regular expression does not permit numeric digits as part of the site name, the pattern matching ends when a digit is reached in the sensor name. In addition to defining

Example 8.11: Calling payload.py

```
<1>$ rwsort inbound.rw --python-file=payload.py --fields=payload \
| rwcut --python-file=payload.py --fields=5,packets,bytes,payload
pro|  packets|    bytes|  payload|
6|    1007|   40280|    0|
1|     1|     28|    0|
6|     2|     92|   12|
6|     8|    332|   12|
1|     1|     48|   20|
17|    1|     51|   23|
1|    14|    784|  392|
80|     3|    762|  702|
50|    16|   1920| 1472|
17|     1|   2982| 2954|
47|    153|  12197| 8525|
50|     77|  11088| 8932|
17|    681| 212153|193085|
6|    309| 398924|386564|
51|   5919| 773077|583669|
51|  10278|1344170|1015274|
6|     820| 1144925|1112125|
97| 2134784|2770949632|2728253952|
```

Example 8.12: decode\_duration.py: A Program to Create a String Field for rwcut

```
'''Define a field for rwcut which formats the flow duration as a string
representation of days, hours, minutes, seconds, and milliseconds.
...

SECPERHOUR = 3600
SECPERMIN = 60

def decode_duration(rec):
    (hours, seconds) = divmod(rec.duration.seconds, SECPERHOUR)
    (minutes, seconds) = divmod(seconds, SECPERMIN)
    return '%02dd%02dh%02dm%02d.%03ds' % (rec.duration.days, hours, minutes,
                                         seconds, rec.duration.microseconds // 1000)

register_field('decode_duration', column_width=16, rec_to_text=decode_duration,
              description='Decomposition of duration into days, hours,
              ' minutes, seconds, and milliseconds')
```

Example 8.13: Calling `decode_duration.py`

```
<1>$ rwcut flows.rw --python-file=decode_duration.py \  
    --fields=decode_duration,duration  
  decode_duration| duration|  
01d07h32m49.161s|113569.161|  
00d00h29m00.450s| 1740.450|  
00d19h05m22.667s|68722.667|  
00d00h00m03.000s|   3.000|
```

a function that derives a string from a SiLK flow record, we need a function to pad the strings so their lengths are the same for all records, and a function to strip the padding.

The functions provided here should work for other string-field plug-ins, except for the call to `get_site`, which derives the string from the record. Then we must establish the maximum length of the field, and supply a few additional parameters to the `register_field` call. Example 8.15 shows how use of the derived field reduces the output to fewer lines than the number of sensors.

---

 Example 8.14: `sitefield.py`: A Program to Create a String Field for Five SiLK Tools
 

---

```

import re # compile(), SRE_Pattern.match(), SRE_Match.group()

# Global variables that may be modified for the enterprise.
MAX_FIELD_LEN = 8
regex = re.compile(r"^[A-Za-z]+")

def get_site(rec):
    '''Derive site name from sensor field in flow record.'''

    return regex.match(rec.sensor).group()

def remove_padding(bin):
    '''Remove padding from fixed-length string.'''

    # for Python 3.x, change "bin" to "str(bin, 'UTF-8')"
    return bin.rstrip('\0')

def pad_to_fixed_length(rec):
    '''Make site names all the same length.'''

    # for Python 3.x, enclose the entire expression in "bytes( , 'UTF-8')"
    return get_site(rec).ljust(MAX_FIELD_LEN, '\0')

register_field('Site',
              bin_bytes=MAX_FIELD_LEN,
              bin_to_text=remove_padding,
              column_width=MAX_FIELD_LEN,
              description='Site name derived from sensor.',
              rec_to_bin=pad_to_fixed_length,
              rec_to_text=get_site)

```

---

Example 8.15: Calling `sitefield.py`

```

<1>$ rwcut flows.rw --python-file=sitefield.py --fields=sensor,Site
  sensor|    Site|
NewYork1| NewYork|
NewYork2| NewYork|
  London0|  London|
London1a|  London|
NewYork1| NewYork|
NewYork2| NewYork|
NewYork3| NewYork|
  London0|  London|
London1a|  London|
London1a|  London|
London1a|  London|
NewYork2| NewYork|
NewYork2| NewYork|
NewYork3| NewYork|
NewYork3| NewYork|
NewYork2| NewYork|
London1a|  London|
<2>$ rwuniq flows.rw --python-file=sitefield.py --field=Site --sort-output
  Site|  Records|
  London|      7|
NewYork|     10|

```

## 8.5 Defining Key Fields and Summary Value Fields for `rwuniq` and `rwstats`

In addition to defining key fields that `rwcut` can use for display, `rwsort` can use for sorting, and `rwgroup` can use for grouping, `rwuniq` and `rwstats` can make use of both key fields and summary value fields. Key fields and summary fields use different registration functions in the simple API, however the registered callback functions do not have to be different. In Example 8.16, the same function, `rec_bpp`, is used to compute the bytes-per-packet ratio for a flow record for use in binning records by a key and in proposing candidate values for the summary value.

Example 8.16: `bpp.py`

```

def rec_bpp(rec):
    return int(round(float(rec.bytes) / float(rec.packets)))

register_int_field('bpp', rec_bpp, 0, (1<<32) - 1)
register_int_max_aggregator('maxbpp', rec_bpp, (1<<32) - 1)

```

In Example 8.17, Command 2 uses the Python file, `bpp.py`, to create a key field for binning records. Command 3 creates an summary value field instead. The summary value in the example finds the maximum value of all the records in a bin, but there are simple API calls for minimum value and sum as well. For additional summaries, the analyst can use the advanced API function, `register_field`.

Example 8.17: Calling bpp.py

```

<1>$ rfilter --type=in --start-date=2015/06/02 \
  --end-date=2015/06/18 --protocol=0- \
  --max-pass-records=70 --pass=tmp.rw
<2>$ rwuniq tmp.rw --python-file=bpp.py --fields=protocol,bpp \
  --values=records
prol      bpp|      Records|
17|        70|         1|
 6|       132|         2|
 6|       410|         1|
17|        74|         2|
 6|       201|         1|
 6|       116|         1|
17|        72|         1|
17|        85|         1|
 6|       409|         1|
 6|       217|         1|
17|       213|         1|
17|        84|         1|
 6|       434|         1|
 6|       140|         1|
 6|        40|        15|
 6|       160|         1|
 6|        46|         1|
 6|       236|         1|
 6|        45|        29|
17|        73|         1|
 6|       174|         5|
 6|       154|         1|
<3>$ rwuniq tmp.rw --python-file=bpp.py --fields=protocol \
  --values=maxbpp
prol      maxbpp|
 6|       434|
17|       213|

```

As shown in this chapter, PySiLK simplifies several previously difficult analyses, without requiring coding large scripts. While the programming involved in creating these scripts has not been described in much detail, the scripts shown (or simple modifications of these scripts) may prove useful to analysts.



# Appendix A

## Networking Primer

This appendix reviews basic topics in Transmission Control Protocol/Internet Protocol (TCP/IP). It is not intended as a comprehensive summary of this topic, but it will help to refresh your knowledge and prepare you for using the SiLK tools for analysis.

Upon completion of this appendix, you will be able to

- describe the structure of IP packets and the relationship between the protocols that constitute the IP protocol suite
- explain the mechanics of TCP, such as the TCP state machine and TCP flags

### A.1 Understanding TCP/IP Network Traffic

This section provides an overview of the TCP/IP networking suite. TCP/IP is the foundation of internetworking. All packets analyzed by the SiLK system use protocols supported by the TCP/IP suite. These protocols behave in a well-defined manner. One possible sign of a security breach is a deviation from accepted behavior. In this section, you will learn about what is specified as accepted behavior. While there are common deviations from the specified behavior, knowing what is specified forms a basis for further knowledge.

This section is intended as a refresher. The TCP/IP suite is a complex collection of more than 50 protocols and comprises far more information than can be covered in this section. A number of online documents and printed books provide other resources on TCP/IP to further your understanding of the TCP/IP suite.

### A.2 TCP/IP Protocol Layers

Figure A.1 shows a basic breakdown of the protocol layers in TCP/IP. The Open Systems Interconnection (OSI) Reference Model, the best known model for layered protocols, consists of seven layers. However, TCP/IP wasn't created with the OSI Reference Model in mind. TCP/IP conforms with the Department of Defense (DoD) ARPANET Reference Model (RFC<sup>17</sup> 871, found at <https://tools.ietf.org/html/rfc871>), a

---

<sup>17</sup>A Request for Comments is an official document, issued by the Internet Engineering Task Force. Some RFCs have Standards status; others do not.

four-layer model. Although TCP/IP and the DoD ARPANET Reference Model have a shared history, it is useful and customary to describe TCP/IP's functions in terms of the OSI Reference Model. OSI is the only model in which network professionals sometimes refer to the layers by number, so any reference to Layer 4, or L4, definitely refers to OSI's Transport layer.

Figure A.1: TCP/IP Protocol Layers

OSI Reference Model	DoD (TCP/IP) Arpanet Ref Model
7 Application	Process Level / Applications
6 Presentation	
5 Session	
4 Transport	Host-to-Host
3 Network	Internet
2 Data-Link	Network
1 Physical	Interface

Starting with the top row of Figure A.1, a *network application* (such as email, telephony, streaming television, or file transfer) creates a *message* that should be understandable by another instance of the network application on another host. This is known as an *application-layer* message. Sometimes the character set, graphics format, or file format must be described to the destination host—as with Multipurpose Internet Mail Extensions (MIME) in email—so the destination host can present the information to the recipient in an understandable way; this is done by adding metadata to the *presentation-layer* header.

Sometimes users want to be able to resume communications sessions when their connections are lost, such as with online games or database updates; this is accomplished with the *session-layer* checkpointing capabilities. Many communications do not use functions of the presentation and session layers, so their headers are omitted. The *transport-layer* protocols identify with port numbers which process or service in the destination host should handle the incoming data; a protocol like User Datagram Protocol (UDP) does little else, but a more complicated protocol like TCP also performs packet sequencing, duplicate packet detection, and lost packet retransmission.

The *network* layer is where we find Internet Protocol, whose job is to route packets from the network interface of the source host to the network interface of the destination host, across many networks and routers in the internetwork. Those networks are of many types (such as Ethernet, Asynchronous Transfer Mode [ATM], cable modem [DOCSIS<sup>®</sup>], or digital subscriber line [DSL]), each with its own frame format and rules described by its *data-link-layer* protocol. The data-link protocol imposes a maximum transmission unit (MTU) size on frames and therefore on datagrams and segments as well. The vast majority of enterprise network traffic is transferred over Ethernet at some point, and Ethernet has the lowest MTU (normally 1,500; 1,492 with

IEEE® 802.2 LLC) of any modern Data-Link layer protocol. So Ethernet’s MTU becomes the effective MTU for the full path.

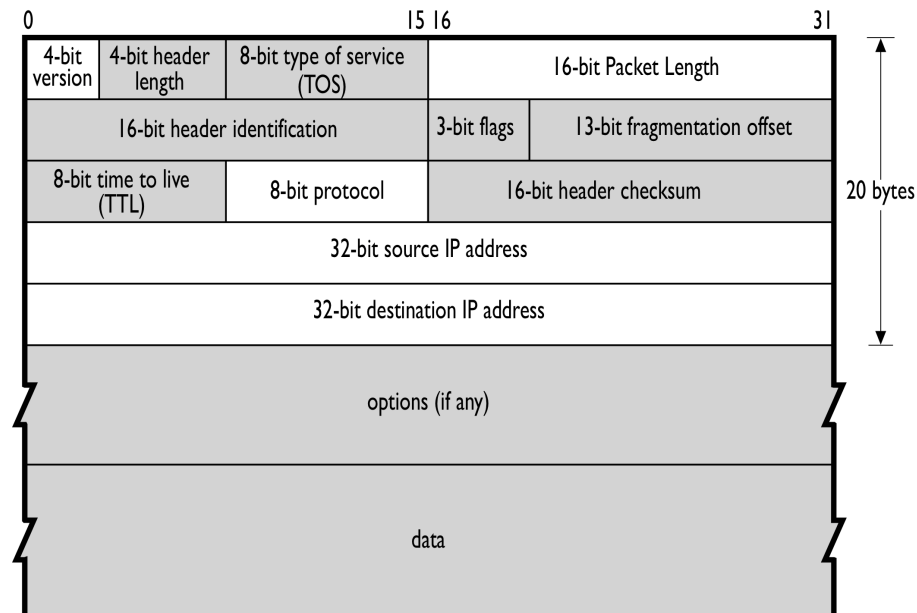
Finally, the frame’s bits are transformed into an energy (electrical, light, or radio wave) signal by the *physical* layer and transmitted across the medium (copper wire, optical fiber, or space). The process of each successively lower layer protocol adding information to the original message is called *encapsulation* because it’s like putting envelopes inside other envelopes.

Each layer adds metadata to the packet that it receives from a higher layer by prepending a header like writing on the outside of that layer’s envelope. When a signal arrives at the destination host’s network interface, the entire process is reversed with *decapsulation*.

### A.3 Structure of the IP Header

IP passes collections of data as datagrams. Two versions of IP are currently used: versions 4 and 6, referred to as IPv4 and IPv6, respectively. IPv4 still constitutes the vast majority of IP traffic in the Internet. IPv6 usage is growing, and both versions are fully supported by the SiLK tools. Figure A.2 shows the breakdown of IPv4 datagrams. Fields that are not recorded by the SiLK data collection tools are grayed out. With IPv6, SiLK records the same information, although the addresses are 128 bits, not 32 bits.

Figure A.2: Structure of the IPv4 Header



## A.4 IP Addressing and Routing

IP can be thought of as a very-high-speed postal service. If someone in Pittsburgh sends a letter to someone in New York, the letter passes through a sequence of postal workers. The postal worker who touches the mail may be different every time a letter is sent, and the only important address is the destination. Normally, there is no reason that New York has to respond to Pittsburgh, and if it does (such as for a return receipt), the sequence of postal workers could be completely different.

IP operates in the same fashion: There is a set of routers between any pair of sites, and packets are sent to the routers the same way that the postal system passes letters back and forth. There is no requirement that the set of routers used to pass data to a destination must be the same as the set used for the return trip, and the routes can change at any time.

Most importantly, the only IP address that must be valid in an IP packet is the destination address. IP itself does not require a valid source address, but some other protocols (e.g., TCP) cannot complete without valid source and destination addresses because the source needs to receive the acknowledgment packets to complete a connection. (However, there are numerous examples of intruders using *incomplete* connections for malicious purposes.)

### A.4.1 Structure of an IP Address

The Internet has space for approximately four billion unique IPv4 addresses. While an IPv4 address can be represented as a 32-bit integer, it is usually displayed in *dotted decimal* (or *dotted quad*) format as a set of four decimal integers separated by periods (dots); for example, 128.2.118.3, where each integer is a number from 0 to 255, representing the value of one byte (octet).

IP addresses and ranges of addresses can also be referenced using *CIDR blocks*. CIDR is a standard for grouping together addresses for routing purposes. When an entity purchases or leases a range of IP addresses from the relevant authorities, that entity buys/leases a routing block, that is used to direct packets to its network.

CIDR blocks are usually described with CIDR notation, consisting of an address, a slash (/), and a prefix length. The prefix length is an integer denoting the number of bits on the left side of the address needed to identify the block. The remaining bits are used to identify hosts within the block. For example, 128.2.0.0/16 would signify that the leftmost 16 bits (2 octets), whose value is 128.2, identify the CIDR block and the remaining bits on the right can have any value denoting a specific host within the block. So all IP addresses from 128.2.0.0 to 128.2.255.255, in which the first 16 bits are unchanged, belong to the same block. Prefix lengths range from 0 (all addresses belong to the same unspecified network; there are 0 network bits specified)<sup>18</sup> to 32 (the whole address is made of unchanging bits, so there is only one address in the block; the address is a single host).

With the introduction of IPv6, all of this is changing. IPv6 addresses are 128 bits in length, for a staggering  $3.4 \times 10^{38}$  (340 undecillion or 340 trillion trillion trillion) possible addresses. IPv6 addresses are represented as groups of eight hexadectets (four hexadecimal digit integers); for example

FEDC:BA98:7654:3210:0037:6698:0000:0510

Each integer is a number between 0 and FFFF (the hexadecimal equivalent of decimal 65,535). IPv6 addresses are allocated in a fashion such that the high-order and low-order digits are manipulated most often, with

<sup>18</sup>CIDR /0 addresses are used almost exclusively for empty routing tables and are not accepted by the SiLK tools. This effectively means the range for CIDR prefix lengths is 1–32 for IPv4.

long strings of hexadecimal zeroes in the middle. There is a shorthand of `::` that can be used once in each address to represent a series of zero-valued hexadectets. The address `FEDC::3210` is therefore equivalent to `FEDC:0:0:0:0:0:0:3210`.

IPv4-compatible (`::0:0/96`) and IPv4-mapped (`::FFFF:0:0/96`) IPv6 addresses are displayed by the SiLK tools in a mixed IPv6/IPv4 format (complying with the canonical format), with the network prefix displayed in hexadecimal, and the 32-bit field containing the embedded IPv4 address displayed in dotted quad decimal. For example, the IPv6 addresses `::102:304` (IPv4-compatible) and `::FFFF:506:708` (IPv4-mapped) will be displayed as `::1.2.3.4` and `::FFFF:5.6.7.8`, respectively.

The routing methods for IPv6 addresses are beyond the scope of this handbook—see RFC 4291 (<https://tools.ietf.org/html/rfc4291>) for a description. Blocks of IPv6 addresses are generally denoted with CIDR notation, just as blocks of IPv4 addresses are. CIDR prefix lengths can range from 0 to 128 in IPv6. For example, `::FFFF:0:0/96` indicates that the most significant 96 bits of the address `::FFFF:0:0` constitute the network prefix (or network address), and the remaining 32 bits constitute the host part.

In SiLK, the support for IPv6 is controlled by configuration. Check for IPv6 support by running `any_SiLK_tool --version` (e.g., `rcut --version`). Then examine the output to see if “IPv6 flow record support” is “yes.”

## A.4.2 Reserved IP Addresses

While IPv4 has approximately four billion addresses available, large segments of IP address space are reserved for the maintenance and upkeep of the Internet. Various authoritative sources provide lists of the segments of IP address space that are reserved. One notable reservation list is maintained by the Internet Assigned Numbers Authority (IANA) at <https://www.iana.org/assignments/ipv4-address-space>. IANA also keeps a list of IPv6 reservations at <https://www.iana.org/assignments/ipv6-address-space>.

In addition to this list, the Internet Engineering Task Force (IETF) maintains several RFCs that specify other reserved spaces. Most of these spaces are listed in RFC 6890, “Special-Purpose IP Address Registries” at <https://tools.ietf.org/html/rfc6890>. Table A.1 summarizes major IPv4 reserved spaces. IPv6 reserved spaces are shown in Table A.2.

Examples in this handbook use addresses in the private and documentation spaces, or addresses that are obviously fictitious, such as 1.2.3.4. This is done to protect the identities of organizations on whose data we tested our examples. Analysts may observe, in real captured traffic, addresses that are not supposed to appear on the Internet. This may be due to misconfiguration of network infrastructure devices or to falsified (spoofed) addressing.

In general, link-local (169.254.0.0/16 in IPv4, FE80::/10 in IPv6) and loopback (127.0.0.0/8 and ::1) destination IP addresses should not cross any routers. Private IP address space (10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16, and FC00::/7) should not enter or traverse the Internet, so it should not appear at edge routers. Consequently, the appearance of these addresses at these routers indicates a failure of routing policy. Similarly, traffic should not come into the enterprise network from these addresses; the Internet as a whole should not route that traffic to the enterprise network.

Table A.1: IPv4 Reserved Addresses

Space	Description	RFC
0.0.0.0/8	This host (0) or specified host on this network (source)	1122
10.0.0.0/8	Private networks	1918
100.64.0.0/10	Carrier-grade Network Address Translation	6598
127.0.0.0/8	Loopback (self-address)	6890
169.254.0.0/16	Link local (autoconfiguration)	6890
172.16.0.0/12	Private networks	1918
192.0.0.0/24	Reserved for IETF protocol assignments	6890
192.0.0.0/29	Dual-Stack Lite	6333
192.0.0.170/31	NAT64/DNS64 Prefix Discovery	7050
192.0.2.0/24	Documentation (example.com or example.net)	5737
192.31.196.0/24	Nameserver for AS112 Redirection Using DNAME	7535
192.52.193.0/24	Automatic Multicast Tunneling	7450
192.88.99.0/24	6to4 relay anycast (border between IPv6 and IPv4)	3068
192.168.0.0/16	Private networks	1918
192.175.48.0/24	Direct Delegation AS112 Service	7534
198.18.0.0/15	Network Interconnect Device Benchmark Testing	2544
198.51.100.0/24	Documentation (example.com or example.net)	5737
203.0.113.0/24	Documentation (example.com or example.net)	5737
224.0.0.0/4	Multicast (destination)	5771
240.0.0.0/4	Future use (except 255.255.255.255)	1112
255.255.255.255	Limited broadcast (destination)	919

Table A.2: IPv6 Reserved Addresses

Space	Description	RFC
::/128	“Unspecified” address (source)	4291
::1/128	Loopback address [similar to 127.0.0.0/8]	4291
::0.0.0.0/96	IPv4-compatible addresses (deprecated by RFC 4291)	1933
::FFFF:0.0.0.0/96	IPv4-mapped addresses	4291
64:FF9B::0.0.0.0/96	IPv4-IPv6 translation with well-known prefix	6052
100::/64	Discard-only address block	6666
2001::/23	IETF protocol assignments	2928
2001::/32	Teredo tunneling	4380
2001:1::1/128	Port Control Protocol Anycast	7723
2001:2::/48	Benchmarking	5180
2001:3::/32	Automatic Multicast Tunneling	7450
2001:4:112::/48	Nameserver for AS112 Redirection Using DNAME	7535
2001:10::/28	Overlay Routable Cryptographic Hash IDentifiers (deprecated)	4843
2001:20::/28	ORCHIDv2	7343
2001:DB8::/32	Documentation addresses [similar to 192.0.2.0/24]	3849
2002::/16	6to4 addresses [related to 192.88.99.0/24]	3056
2620:4F:8000::/48	Direct Delegation AS112 Service	7534
FC00::/7	Unique local addresses [similar to RFC 1918 private addresses] primarily seen as FD00::/8	4193
FE80::/10	Link-local unicast (similar to 169.254.0.0/16)	4291
FEC0::/10	Formerly reserved for site-local unicast addresses (deprecated by RFC 3879)	1884
FF00::/8	Multicast [similar to 224.0.0.0/4]	4291

## A.5 Major Protocols

### A.5.1 Protocol Layers and Encapsulation

In the multi-layered scheme used by TCP/IP, lower layer protocols *encapsulate* higher layer protocols, like envelopes within envelopes. When we open the innermost envelope, we find the message that belongs to the highest layer protocol. Conceptually, the envelopes have metadata written on them. In practice, the metadata are recorded in *headers*. The header for the lowest layer protocol is sent over the network first, followed by the headers for progressively higher layers. Finally, the message from the highest layer protocol is sent after the last header.

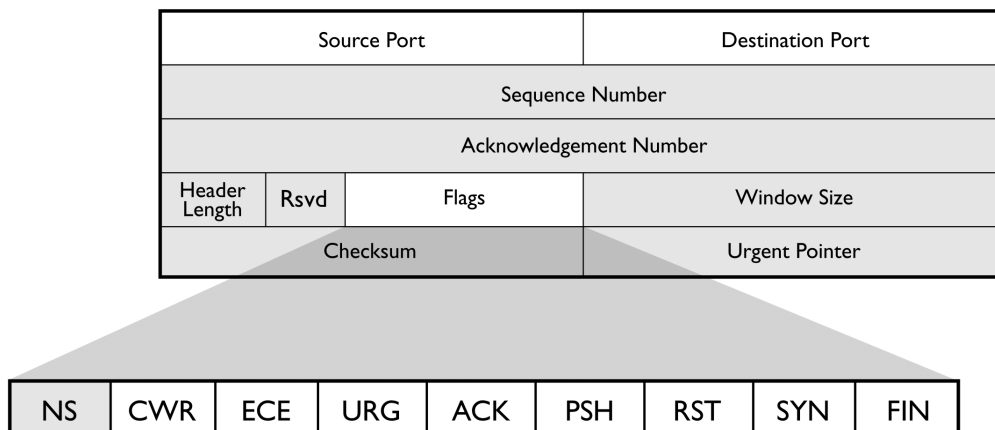
TCP/IP was created before the OSI Reference Model. But if we refer to a layer by its number (e.g., Layer 3 or L3), we always mean the specified layer in that model. While the preceding description of encapsulation is generally true, the model actually assigns protocols to layers based on the protocol's functions, not its order of encapsulation. This is most apparent with Internet Control Message Protocol (ICMP), which the model assigns to the Network layer (L3), even though its header and payload are encapsulated by IP, which is also a Network layer protocol. From here on, we will ignore this fine distinction, and we will consider ICMP to be a Transport layer (L4) protocol because it is encapsulated by IP, a Layer 3 protocol.

### A.5.2 Transmission Control Protocol (TCP)

TCP, the most commonly encountered transport protocol on the Internet, is a stream-based protocol that reliably transmits data from the source to the destination. To maintain this reliability, TCP is very complex: The protocol is slow and requires a large commitment of resources.

Figure A.3 shows a breakdown of the TCP header, which adds 20 additional bytes to the IP header. Consequently, TCP packets will always be at least 40 bytes (60 for IPv6) long. As the shaded portions of Figure A.3 show, most of the TCP header information is not retained in SiLK flow records.

Figure A.3: TCP Header





TCP is built on top of an unreliable infrastructure provided by IP. IP assumes that packets can be lost without a problem, and that responsibility for managing packet loss is incumbent on services at higher layers. TCP, which provides ordered and reliable streams on top of this unreliable packet-passing model, implements this feature through a complex state machine as shown in Figure A.4. The transitions in this state machine are described by labels in a  $\frac{\text{stimulus}}{\text{action}}$  format, where the top value is the stimulating event and the bottom values are actions taken prior to entry into the destination state. Where no action takes place, an “x” is used to indicate explicit inaction.

This handbook does not thoroughly describe the state machine in Figure A.4 (see <https://tools.ietf.org/html/rfc793> for a complete description), however, flows representing well-behaved TCP sessions will behave in certain ways. For example, a flow for a complete TCP session must have at least four packets: one packet that sets up the connection, one packet that contains the data, one packet that terminates the session, and one packet acknowledging the other side’s termination of the session.<sup>19</sup> TCP behavior that deviates from this provides indicators that can be used by an analyst. An intruder may send packets with odd TCP flag combinations as part of a scan (e.g., with all flags set on). Different operating systems handle protocol violations differently, so odd packets can be used to elicit information that identifies the operating system in use or to pass through some systems benignly, while causing mischief in others.

## TCP Flags

TCP uses *flags* to transmit state information among participants. A flag has two states: high or low; so a flag represents one bit of information. There are six commonly used flags:

**ACK:** Short for “acknowledge,” ACK flags are sent in almost all TCP packets and used to indicate that previously sent packets have been received.

**FIN:** Short for “finalize,” the FIN flag is used to terminate a session. When a packet with the FIN flag is sent, the target of the FIN flag knows to expect no more input data. When both have sent and acknowledged FIN flags, the TCP connection is closed gracefully.

**PSH:** Short for “push,” the PSH flag is used to inform a TCP receiver that the data sent in the packet should immediately be sent to the target application (i.e., the sender has completed this particular send), approximating a message boundary in the stream.

**RST:** Short for “reset,” the RST flag is sent to indicate that a session is incorrect and should be terminated. When a target receives a RST flag, it terminates immediately. Some implementations terminate sessions using RST instead of the more proper FIN sequence.

**SYN:** Short for “synchronize,” the SYN flag is sent at the beginning of a session to establish initial sequence numbers. Each side sends one SYN packet at the beginning of a session.

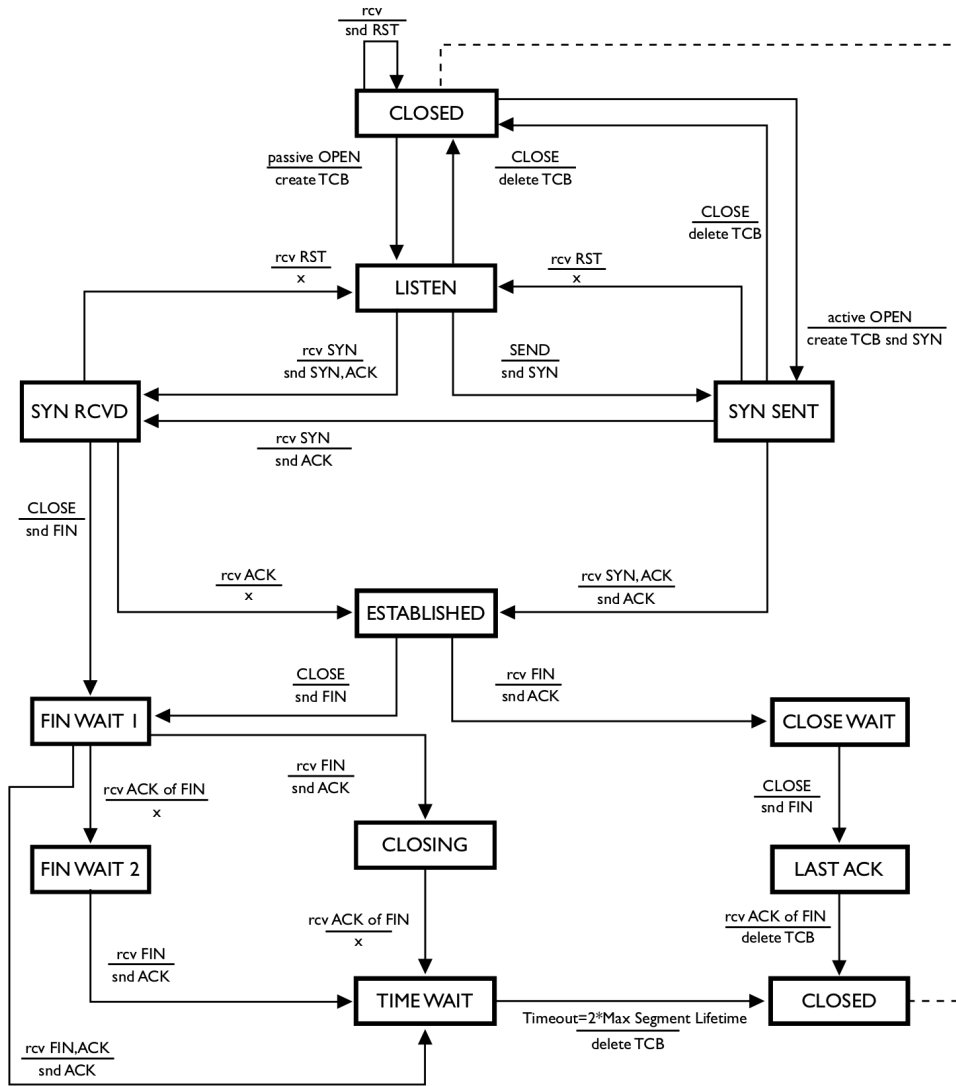
**URG:** Short for “urgent” data, the URG flag is used to indicate that urgent data (such as a signal from the sending application) is in the buffer and should be used first. The URG flag should only be seen in Telnet-like protocols such as Secure Shell (SSH). Tricks with URG flags can be used to fool intrusion detection systems (IDS).

Reviewing the state machine will show that most state transitions are handled through the use of SYN, ACK, FIN, and RST. The PSH and URG flags are less directly relevant. Two other rarely used flags are understood

---

<sup>19</sup>It is technically possible for there to be a valid three-packet complete TCP flow: one SYN packet, one SYN-ACK packet containing the data, and one RST packet terminating the flow. This is a very rare circumstance; most complete TCP flows have more than four packets.

Figure A.4: TCP State Machine



by SiLK: ECE (Explicit Congestion Notification Echo) and CWR (Congestion Window Reduced). Neither is relevant to security analysis at this time, although they can be used with the SiLK tool suite if required. A ninth TCP flag, NS (Nonce Sum), is not recognized or supported by SiLK.

## Major TCP Services

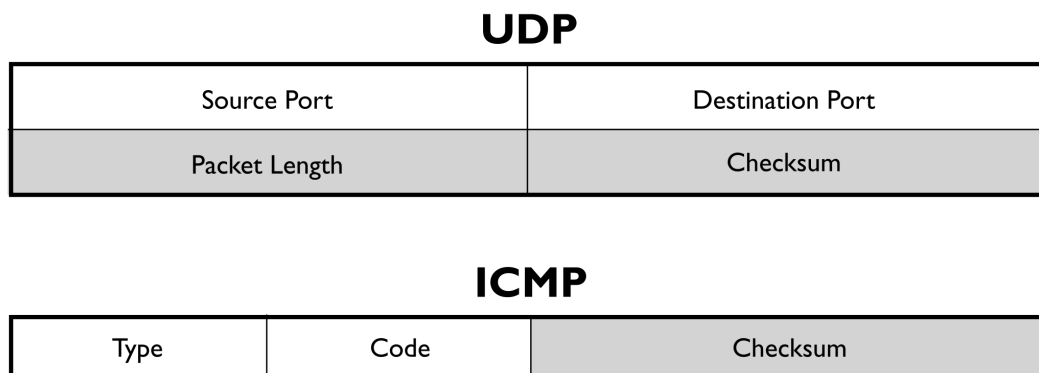
Traditional TCP services have well-known ports; for example, 80 is Web, 25 is SMTP, and 53 is DNS. IANA maintains a list of these port numbers at <https://www.iana.org/assignments/service-names-port-numbers>. This list is useful for legitimate services, but it does not necessarily contain new services or accurate port assignments for rapidly changing services such as those implemented via peer-to-peer networks. Furthermore, there is no guarantee that traffic seen (e.g., on port 80) is actually web traffic or that web traffic cannot be sent on other ports.

### A.5.3 UDP and ICMP

After TCP, the most common protocols on the Internet are UDP and ICMP. While IP uses its addressing and routing to deliver packets to the correct interface on the correct host, Transport layer protocols like TCP and UDP use their port numbers to deliver packets inside the host to the correct process or service. Whereas TCP also provides other functions, such as data streams and reliability, UDP provides only delivery. UDP does not understand that sequential packets might be related (as in streams); UDP leaves that up to higher layer protocols. UDP does not provide reliability functions, like detecting and recovering lost packets, reordering packets, or eliminating duplicate packets. UDP is a fast but unreliable message-passing mechanism used for services where throughput is more critical than accuracy. Examples include audio/video streaming, as well as heavy-use services such as the Domain Name System (DNS).

ICMP, a reporting protocol that works in tandem with IP, sends error messages and status updates, and provides diagnostic capabilities like echo.

Figure A.5: UDP and ICMP Headers



## UDP and ICMP Packet Structure

Figure A.5 shows a breakdown of UDP and ICMP packets, as well as the fields collected by SiLK. UDP can be thought of as TCP without the additional state mechanisms; a UDP packet has both source and destination ports, assigned in the same way TCP assigns them, as well as a payload.

ICMP is a straight message-passing protocol and includes a large amount of information in its first two fields: Type and Code. The Type field is a single byte indicating a general class of message, such as “destination unreachable.” The Code field contains a byte indicating greater detail about the type, such as “port unreachable.” ICMP messages generally have a limited payload; most messages have a fixed size based on type, with the notable exceptions being echo request (ICMPv4 type 8 or ICMPv6 type 128) and echo reply (ICMPv4 type 0 or ICMPv6 type 129).

## Major UDP Services and ICMP Messages

UDP services are covered in the IANA webpage whose URL is listed above. As with TCP, the values given by IANA are slightly behind those currently observed on the Internet. IANA also excludes port utilization (even if common) by malicious software such as worms. Although not official, numerous port databases on the web can provide insight into the current port utilization by services.

ICMPv4 types and codes are listed at <https://www.iana.org/assignments/icmp-parameters>. ICMPv6 types and codes are listed at <https://www.iana.org/assignments/icmpv6-parameters>. These lists are definitive and include references to RFCs explaining the types and codes.

## Appendix B

# Using UNIX to Implement Network Traffic Analysis

This appendix provides a review of basic UNIX operations. SiLK is implemented on UNIX (e.g., Apple® OS X®, FreeBSD®, Solaris®) and UNIX-like operating systems and environments (e.g., Linux®, Cygwin); consequently an analyst must be able to work with UNIX to use the SiLK tools.

### B.1 Using the UNIX Command Line

UNIX uses a program known as a shell to obtain commands from a user and either perform the task described by that command or invoke another program that will. Linux usually uses Bash (Bourne-Again SHell) for its shell. When the shell is ready to accept a command from the user, it displays a string of characters known as a prompt to let the user know that he or she can enter a command now. Besides notifying the user that a command can be accepted at this time, the prompt may convey additional information. The choice of information to be conveyed may be made by the user by providing a prompt template to the shell. In this handbook, the prompt will appear as in Example B.1.

Example B.1: A UNIX Command Prompt

---

```
<1>$
```

---

The integer between angle brackets will be used to refer to specific commands in examples. Commands can be invoked by typing them directly at the command line. UNIX commands are typically abbreviated English words and accept space-separated parameters. Parameters are just values (like filenames), an option-name/value pair, or just an option name. Option names are double dashes followed by hyphenated words, single dashes followed by single letters, or (rarely) single dashes followed by words (as in the `find` command). Table B.1 lists some of the more common UNIX commands. To see more information on these commands type `man` followed by the command name. Example B.2 and the rest of the examples in this handbook show the use of some of these commands.

Table B.1: Some Common UNIX Commands

Command	Description
<code>cat</code>	Copies streams and/or files onto standard output (show file content)
<code>cd</code>	Changes [working] directory
<code>chmod</code>	Changes file-access permissions. Needed to make script executable
<code>cp</code>	Copies a file from one name or directory to another
<code>cut</code>	Isolates one or more columns from a file
<code>date</code>	Shows the current or calculated day and time
<code>echo</code>	Writes arguments to standard output
<code>exit</code>	Terminates the current shell or script (log out) with an exit code
<code>export</code>	Assigns a value to an environment variable that programs can use
<code>file</code>	Identifies the type of content in a file
<code>grep</code>	Displays from a file those lines matching a given pattern
<code>head</code>	Shows the first few lines of a file's content
<code>kill</code>	Terminates a job or process
<code>less</code>	Displays a file one full screen at a time
<code>ls</code>	Lists files in the current (or specified) directory -l (for long) parameter to show all directory information
<code>man</code>	Shows the online documentation for a command or file
<code>mkdir</code>	Makes a directory
<code>mv</code>	Renames a file or moves it from one directory to another
<code>ps</code>	Displays the current processes
<code>pwd</code>	Displays the working directory
<code>rm</code>	Removes a file
<code>sed</code>	Edits the lines on standard input and writes them to standard output
<code>sort</code>	Sorts the contents of a text file into lexicographic order
<code>tail</code>	Shows the last few lines of a file
<code>time</code>	Shows the execution time of a command
<code>top</code>	Shows the running processes with the highest CPU utilization
<code>uniq</code>	Reports or omits repeated lines. Optionally counts repetitions
<code>wc</code>	Counts the words (or, with -l parameter, counts the lines) in a file
<code>which</code>	Verifies which copy of a command's executable file is used
<code>\$(...)</code>	Inserts the output of the contained command into the command line
<code>var=value</code>	Assigns a value to a shell variable. For use by the shell only, not programs

## Example B.2: Using Simple UNIX Commands

```
<1>$ echo Here are some simple commands:
Here are some simple commands:
<2>$ date
Thu Jul  3 15:56:24 EDT 2014
<3>$ date -u
Thu Jul  3 19:56:24 UTC 2014
<4>$ # This is a comment line. It has no effect.
<5>$ #The next command lists my running processes
<6>$ ps -f
UID          PID    PPID  C STIME TTY          TIME CMD
user1         8280   8279  0 14:43 pts/2        00:00:00 -bash
user1        10358  10355  1 15:56 pts/2        00:00:00 ps -f
<7>$ cat animals.txt
Animal  Legs   Color
-----  ----  -----
fox      4      red
gorilla  2      silver
spider   8      black
moth     6      white
<8>$ file animals.txt
animals.txt: ASCII text
<9>$ head -n 3 animals.txt
Animal  Legs   Color
-----  ----  -----
fox      4      red
<10>$ cut -f 1,3 animals.txt
Animal  Color
-----  -----
fox     red
gorilla silver
spider  black
moth    white
```

## B.2 Standard In, Out, and Error

Many UNIX programs, including most of the SiLK tools, have a default for where to obtain input and where to write output. The symbolic filenames `stdin`, `stdout`, and `stderr` are not the names of disk files, but rather they indirectly refer to files. Initially, the shell assigns the keyboard to `stdin` and assigns the screen to `stdout` and `stderr`. Programs that were written to read and write through these symbolic filenames will default to reading from the keyboard and writing to the screen. But the symbolic filenames can be made to refer indirectly to other files, such as disk files, through shell features called *redirection* and *pipes*.

### B.2.1 Output Redirection

Some programs, like `cat` and `cut`, have no way for the user to tell the program *directly* which file to use for output. Instead these programs always write their output to `stdout`. The user must inform *UNIX*, not the program, that `stdout` should refer to the desired file. The program then only knows its output is going to `stdout`, and it's up to UNIX to route the output to the desired file. One effect of this is that any error message emitted by the program that refers to its output file can only display “`stdout`,” since the actual output filename is unknown to the program.

The shell makes it easy to tell UNIX that you wish to *redirect* `stdout` from its default (the screen) to the file that the user specifies. This is done right on the same command line that runs the program, using the greater-than symbol (`>`) and the desired filename (as shown in Command 1 of Example B.3).

SiLK tools that write binary (non-text) data to `stdout` will emit an error message and terminate if `stdout` is assigned to a terminal device. Such tools must have their output directed to a disk file or piped to a SiLK tool that reads that type of binary input.

Example B.3: Output Redirection

---

```

<1>$ cut -f 1,3 animals.txt >animalcolors.txt
<2>$ cat animalcolors.txt
Animal  Color
-----  -----
fox     red
gorilla silver
spider  black
moth   white
<3>$ rm animalcolors.txt
<4>$ ls animalcolors.txt
ls: animalcolors.txt: No such file or directory

```

---

### B.2.2 Input Redirection

A very few programs, like `tr`, have no syntax for specifying the input file and rely entirely on UNIX to connect an input file to `stdin`. The shell provides a method for redirecting input very similar to redirecting output. You specify a less-than symbol (`<`) followed by the input filename as shown in Command 2 of Example B.4.



## Example B.4: Input Redirection

---

```
<1>$ #Translate hyphens to slashes
<2>$ tr - / <animals.txt
Animal  Legs   Color
//////// //    //
fox     4     red
gorilla 2     silver
spider  8     black
moth   6     white
```

---

## B.2.3 Pipes

The real power of `stdin` and `stdout` becomes apparent with *pipes*. A pipe connects the `stdout` of the first program to the `stdin` of a second program. This is specified in the shell using a vertical bar character (`|`), known in UNIX as the pipe symbol.

## Example B.5: Using a Pipe

---

```
<1>$ head -n 4 animals.txt | cut -f 1,3
Animal  Color
-----  -----
fox     red
gorilla silver
```

---

In Example B.5, the `head` program reads the first four lines from the `animals.txt` file and writes those lines to `stdout` as normal, except that `stdout` does not refer to the screen. The `cut` program has no input filename specified and is programmed to read from `stdin` when no input filename appears on the command line. The pipe connects the `stdout` of `head` to the `stdin` of `cut` so that `head`'s output lines become `cut`'s input lines without those lines ever touching a disk file. `cut`'s `stdout` was not redirected, so its output appears on the screen.

## B.2.4 Here-Documents

Sometimes we have a small set of data that is manually edited and perhaps doesn't change from one run of a script to the next. If so, instead of creating a separate data file for the input, we can put the input data right into the script file. This is called a *here-document*, because the data are right here in the script file, immediately following the command that reads them.

Example B.6 illustrates the use of a here-document to supply several filenames to a SiLK program called `rwsort`. The `rwsort` program has an option called `--xargs` telling it to get a list of input files from `stdin`. The here-document supplies data to `stdin` and is specified with double less-than symbols (`<<`), followed by a string that defines the marker that will indicate the end of the here-document data. The lines of the script file that follow the command are input data lines until a line with the marker string is reached.

## Example B.6: Using a Here-Document

---

```
<1>$ rwsort --xargs --fields=sTime --output-path=week.rw <<END-OF-LIST
sunday.rw
monday.rw
tuesday.rw
wednesday.rw
thursday.rw
friday.rw
saturday.rw
END-OF-LIST
<2>$ rfileinfo --fields=count-records *day.rw week.rw
```

---

## B.2.5 Named Pipes

Using the pipe symbol, a script creates an *unnamed* pipe. Only one unnamed pipe can exist for output from a program, and only one can exist for input to a program. For there to be more than one, you need some way to distinguish one from another. The solution is *named* pipes.

Unlike unnamed pipes, which are created in the same command line that uses them, named pipes must be created prior to the command line that employs them. As named pipes are also known as *FIFOs* (for First In First Out), the command to create one is `mkfifo` (make FIFO). Once the FIFO is created, it can be opened by one process for reading and by another process (or multiple processes) for writing.

Scripts that use named pipes often employ another useful feature of the shell: running programs in the background. In Bash, this is specified by appending an ampersand (&) to the command line. When a program runs in the background, the shell will not wait for its completion before giving you a command prompt. This allows you to issue another command to run concurrently with the background program. You can force a script to wait for the completion of background programs before proceeding by using the `wait` command.

SiLK applications can communicate via named pipes. In Example B.7, we create a named pipe (in Command 1) that one call to `rwfilter` (in Command 2) uses to filter data concurrently with another call to `rwfilter` (in Command 3). Results of these calls are shown in Commands 5 and 6. Using named pipes, sophisticated SiLK operations can be built in parallel. A backslash (\) at the very end of a line indicates that the command is continued on the following physical line.

## Example B.7: Using a Named Pipe

```

<1>$ mkfifo /tmp/namedpipe1
<2>$ rfilter --start-date=2014/03/21T17 --end-date=2014/03/21T18 \
           --type=all --protocol=6 \
           --fail=/tmp/namedpipe1 --pass=stdout \
           | rwuniq --fields=protocol --output-path=tcp.out &
<3>$ rfilter /tmp/namedpipe1 --protocol=17 --pass=stdout \
           | rwuniq --fields=protocol --output-path=udp.out &
<4>$ wait
<5>$ cat tcp.out
prol  Records|
  6|  34866860|
<6>$ cat udp.out
prol  Records|
 17|  17427015|
<7>$ rm /tmp/namedpipe1 tcp.out udp.out

```

## B.3 Script Control Structures

Some advanced examples in this handbook will use control structures available from Bash. The syntax

```
for name in word-list-expression; do ... done
```

indicates a loop where each of the space-separated values returned by *word-list-expression* is given in turn to the variable indicated by *name* (and referenced in commands as *\$name*), and the commands between *do* and *done* are executed with that value. The syntax

```
while expression; do ... done
```

indicates a loop where the commands between *do* and *done* are executed as long as *expression* evaluates to true.



# Appendix C

## SiLK Commands

This appendix lists the SiLK commands that are used in this guide and describes their most commonly-used options. Section [C.27](#) surveys features, including parameters, that are common across several tools.

### C.1 Getting Help with SiLK Tools

All SiLK tools include a help screen that provides a summary of command information. To view the help screen, specify the `--help` parameter with the command.

```
$ command --help
```

SiLK is distributed with UNIX manual pages that explain all the parameters and functionality of each tool in the suite.

```
$ man command
```

All SiLK tools also have a `--version` parameter that identifies the version installed. Since the suite is still being extended and evolved, this version information may be quite important.

```
$ command --version
```

## C.2 rwsiteinfo Command Summary

### rwsiteinfo

<b>Description</b>	Displays SiLK configuration file information for a site, its sensors, and the traffic they collect.
<b>Call</b>	<code>rwsiteinfo --fields=<i>parameters</i> --sensor=<i>sensors</i></code>
<b>Parameters</b>	<p><b>--fields</b> Displays information about the parameters specified in a comma-separated list (required). See Table C.1 for a list of commonly-used parameters. Specify <code>:list</code> after a parameter to display output in a comma-separated list instead of columns.</p> <p><b>--sensor</b> Displays information about the specified sensor or sensors</p> <p><b>--type</b> Displays information about the SiLK types specified in a comma-separated list</p> <p><b>--column-separator</b> Uses the specified character as a column separator (default is  )</p> <p><b>--no-titles</b> Do not print column headers.</p> <p><b>--list-delimiter</b> Uses the specified character as a list delimiter (default is ,)</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

Table C.1: Parameters for `rwsiteinfo --fields`

Parameter	Description
<code>sensor</code>	Name of sensor
<code>describe-sensor</code>	Description of sensor from configuration file
<code>type</code>	Type of SiLK network data (see Section 1.2.6 for more information)
<code>repo-start-date</code>	Time and date of first data file written to SiLK repository
<code>repo-send-date</code>	Time and date of last data file written to SiLK repository
<code>repo-file-count</code>	Number of files written to SiLK repository

## C.3 rfilter Command Summary

### rfilter

<b>Description</b>	Searches for, retrieves, and partitions SiLK flow records
<b>Call</b>	<code>rfilter {selection   input} partition output [other]</code>
<b>Parameters</b>	<p>Specify either <i>selection</i> or <i>input</i> parameters.</p> <p><i>selection</i> parameters (described in Table C.2) tell <code>rfilter</code> to pull data from SiLK repository files by supplying desired attributes of the records stored in the repository.</p> <p><i>input</i> parameters specify a SiLK data source other than the repository and can include <i>filenames</i> (e.g., <code>infile.rw</code>) or <i>pipe names</i> (e.g., <code>stdin</code> or <code>/tmp/my.fifo</code>) to specify locations from which to read records. As many filenames as desired may be given, with both files and pipes used in the same command. The <code>--xargs</code> parameter specifies a file containing filenames from which to read flow records.</p> <p><i>partition</i> parameters define tests that divide the input records into two groups: (1) “pass” records, which satisfy all the tests and (2) “fail” records, which fail to satisfy at least one of the tests. Each call to <code>rfilter</code> must have at least one partitioning parameter unless the only output parameter is <code>--all-destination</code>, which does not require the difference between passing and failing to be defined. Commonly used partitioning parameters are listed in Tables C.3–C.9.)</p> <p><i>output</i> specify which statistics or sets of records should be returned from the call. There are five output parameters, as described in Table C.10. Each call to <code>rfilter</code> must have at least one of these parameters.</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

#### rfilter Selection Parameters

Table C.2: rfilter Selection Parameters

Parameter	Example	Description
<code>--data-rootdir</code>	<code>/datavol/repos</code>	Location of SiLK repository
<code>--sensors*</code>	<code>1-5</code>	Sensor(s) used to collect data
<code>--class*</code>	<code>all</code>	Category of sensors
<code>--type*</code>	<code>inweb,in,outweb,out</code>	Category of flows within class
<code>--flowtypes*</code>	<code>c1/in,c2/all</code>	Class/type pairs
<code>--start-date</code>	<code>2014/6/13</code>	First day or hour of data to examine
<code>--end-date</code>	<code>2014/3/20T23</code>	Final day or hour of data to examine

\*These selection keywords may be used as partitioning options if an input file or pipe is named.

## rwfilter Partitioning Parameters

Figure C.1 shows several groups of partitioning parameters. This section focuses on the parameters that partition records based on fields in the flow records. Section 5.1.3 discusses IP sets and how to filter with those sets. Section 6.2.7 describes prefix maps and country codes. Section 6.2.1 discusses tuple files and the parameters that use them. Lastly, Section 8.1.2 describes the use of PySiLK plug-ins. Figure C.1 also illustrates the relative efficiency of the types of partitioning parameters. Choices higher in the illustration may be more efficient than choices lower down. So analysts should prefer IPsets to tuple files when either would work, and they should prefer `--address` or `--scidr` to IPsets when those would all work.

Figure C.1: `rwfilter` Partitioning Parameters

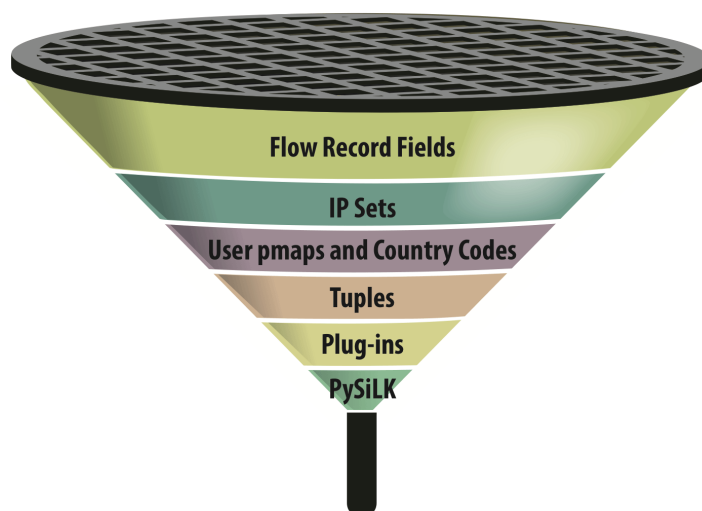


Table C.3: Single-Integer- or Range-Partitioning Parameters

Parameter	Example	Partition Based on
<code>--packets</code>	1–3	Packet count in flow
<code>--bytes</code>	400–2400	Byte count in flow
<code>--bytes-per-packet</code>	1000–1400	Average bytes/packet in flow
<code>--duration</code>	1200–	Duration of flow in seconds
<code>--ip-version</code>	6	IP version of the flow record



Table C.4: Multiple-Integer- or Range-Partitioning Parameters

Parameter	Example	Partition Based on
--protocol	0,2-5,7-16,18-	Protocol number (1=ICMP, 6=TCP, 17=UDP)
--sport	0-1023	Source port
--dport	25	Destination port
--aport	80,8080	Any port. Like --sport, but for either source or destination
--application	2427,2944	Application-layer protocol (see <a href="https://tools.netsa.cert.org/yaf/applabel.html#LABELS">https://tools.netsa.cert.org/yaf/applabel.html#LABELS</a> )
--icmp-type	0-41,253,254	Type of ICMP message
--icmp-code	0,16	Subtype of ICMP message

Table C.5: Address-Partitioning Parameters

Parameter	Example	Partition Based on
--saddress	198.51.100.1,254	Single address, CIDR block, or wildcard for source
--daddress	198.51.100.0/24	Like --saddress, but for destination
--any-address	2001:DB8::x	Like --saddress, but for either source or destination
--next-hop-id	10.2-5.x.x	Like --saddress, but for next hop address
--scidr	198.51.100.1,203.0.113.64/29	Multiple addresses and CIDR blocks for source address
--dcidr	FC00::/7,2001:DB8::/32	Like --scidr, but for destination
--any-cidr	203.0.113.199,192.0.2.44	Like --scidr, but for either source or destination
--nhcidr	203.0.113.8/30,192.0.2.6	Like --scidr, but for next hop address
--sipset	tornodes.set	Source address existing in IPset file
--dipset	websvrs.set	Destination address existing in IPset file
--anyset	dnssvrs.set	Either source or destination address in IPset file
--nhipset	lgvolume.set	Next hop address in IPset file
--not-param	Any address parameter can be inverted using the --not- prefix	

Table C.6: High/Mask Partitioning Parameters

Parameter	Example	Partition Based on
--flags-all	SF/SF,SR/SR	Accumulated TCP flags
--flags-initial	S/SA	TCP flags in the first packet of flow
--flags-session	/FR	Flags in the packets after the first
--attributes	T/T,C/C	Termination attributes or packet size uniformity

Table C.7: Time-Partitioning Parameters

Parameter	Example	Partition Based on
<code>--stime</code>	2014/4/7-2014/4/8T12	Flow's start time
<code>--etime</code>	2014/4/7T8:30-2014/4/8T8:59	Flow's end time
<code>--active-time</code>	2014/4/7T11:33:30-	Overlap of active range and period between flow's start and end times

Table C.8: Prefix-Map-Partitioning Parameters

Parameter	Example	Partition Based on
<code>--pmap-src-mapname</code>	Zer0n3t,YOLOBotnet	Source mapping to a specified label
<code>--pmap-dst-mapname</code>	DMZ,bizpartner	Destination mapping to a specified label
<code>--pmap-any-mapname</code>	mynetdvcs	Source or destination mapping to a specified label (see Section 6.2.7 on page 121)
<code>--scc</code>	ru,cn,br,ko	Source address's country code
<code>--dcc</code>	ca,us,mx	Destination address's country code
<code>--any-cc</code>	a1,a2,o1,--	Source or destination address's country code
<code>--stype</code>	1	category index of source address
<code>--dtype</code>	2	category index of destination address

Table C.9: Miscellaneous Partitioning Parameters

Parameter	Example	Partition Based on
<code>--tuple-file</code>	torguard.tuple	Match to any specified combination of field values (see Section 6.2.1 on page 107)
<code>--python-expr</code>	'rec.sport==rec.dport'	Truth of expression (see Section 8.2 on page 143)
<code>--python-file</code>	complexrules.py	Truth of value returned by the program in the Python file (see Section 8.2 on page 143)
<code>--plugin</code>	flowrate.so	Custom partitioning with a C language program (see Section 8.1 on page 142)

**Note on specifying partitioning parameters:** Partitioning parameters specify a collection of flow record criteria, such as the protocols 6 and 17 or the specific IP address 198.51.100.71. As a result, almost all partitioning parameters describe some group of values. These ranges are generally expressed in the following ways:

**Value range:** Value ranges are used when all values in a closed interval are desired. A value range is two numbers separated by a hyphen, such as `--packets=1-4`, which indicates that flow records with a packet count from one through four (inclusive) belong to the *pass* set of records. Some partitioning parameters (such as `--duration`) only make sense with a value range (searching for flows with a duration exact to the millisecond would be fruitless); An omitted value on the end of the range (e.g., `--bytes=2048-`) specifies that any value greater than or equal to the low value of the range passes. Omitting the value at the start of a range is not permitted. The options in Table C.3 should be specified with a single value range, except `--ip-version` which accepts a single integer.

**Value alternatives:** Fields that have a finite set of values (such as ports or protocol) can be expressed using a comma-separated list. In this format a field is expressed as a set of numbers separated by commas. When only one value is acceptable, it is presented without a comma. Examples include `--protocol=3` and `--protocol=3,9,12`. Value ranges can be used as elements of value alternative lists. For example, `--protocol=0,2-5,7-16,18-` says that all flow records that are not for ICMP (1), TCP (6), or UDP (17) traffic are desired. The options in Table C.4 are specified with value alternatives.

**IP addresses:** IP address specifications are expressed in three ways: a single address; a CIDR block (a network address, a slash (/), and a prefix length); and a SiLK address wildcard. The `--saddress`, `--daddress`, `--any-address`, and `--next-hop-id` options, and the `--not-` forms of those options accept a single specification that may be any of three forms. The `--scidr`, `--dcidr`, `--any-cidr`, and `--nhcidr` options, and the `--not-` forms of those options accept a comma-separated list of specifications that may be addresses or CIDR blocks but not wildcards. SiLK address wildcards are address specifications with a syntax unique to the SiLK tool suite. Like CIDR blocks, a wildcard specifies multiple IP addresses. But while CIDR blocks only specify a continuous range of IP addresses, a wildcard can even specify discontinuous addresses. For example, the wildcard `1-13.1.1.1,254` would select the addresses 1.1.1.1, 2.1.1.1, and so on until 13.1.1.1, as well as 1.1.1.254, 2.1.1.254, and so on until 13.1.1.254. For convenience, the letter *x* can be used to indicate all values in a section (equivalent to `0-255` in IPv4 addresses, `0-FFFF` in IPv6 addresses). CIDR notation may also be used, so `1.1.0.0/16` is equivalent to `1.1.x.x` and `1.1.0-255.0-255`. Any address range that can be expressed with CIDR notation also can be expressed with a wildcard. Since CIDR notation is more widely understood, it probably should be preferred for those cases. As explained in Section A.4.1, IPv6 addresses use a double-colon syntax as a shorthand for any sequence of zero values in the address and use a CIDR prefix length. The options in Table C.5, except the set-related options, are specified with IP addresses. The set-related options specify a filename.

**TCP flags:** The `--flags-all`, `--flags-initial`, and `--flags-session` options to *rwfilter* use a compact, yet powerful, way of specifying filter predicates based on the presence of TCP flags. The argument to this parameter has two sets of TCP flags separated by a forward slash (/). To the left of the slash is the *high* set; it lists the flags that must be set for the flow record to pass the filter. The flag set to the right of the slash contains the *mask*; this set lists the flags whose status is of interest, and the set must be non-empty. Flags listed in the mask set but not in the high set must be off to pass. The flags listed in the high set must be present in the mask set. (For example, `--flags-initial=S/SA` specifies a filter for flow records that initiate a TCP session; the S flag is high [on] and the A flag is low [off].) The options in Table C.6, except `--attributes`, are specified with TCP flags; `--attributes` also specifies flags in a high/mask format, but the flags aren't TCP flags.

**Attributes:** The `--attributes` parameter takes any combination of the letters **S**, **T**, and **C**, expressed in high/mask notation just as for TCP flags. **S** indicates that all packets in the flow have the same length (never present for single-packet flows). **T** indicates the collector terminated the flow collection due to active timeout. **C** indicates the collector produced the flow record to continue flow collection that was terminated due to active timeout. Only the `--attributes` parameter uses attributes for values.

**Time ranges:** Time ranges are two times, potentially precise to the millisecond, separated by a hyphen; in SiLK, these times can be expressed in their full `YYYY/MM/DDThh:mm:ss.mmm` form (e.g., `2005/02/11T03:18:00.005-2005/02/11T05:00:00.243`). The times in a range may be abbreviated by omitting a time (but not date) component and all the succeeding components. If all the time components are omitted, the **T** (or colon) that separates the time from the date may also be omitted. Abbreviated times (times without all the components down to the millisecond) are treated as though the last component supplied includes the entire period specified by that component, not just an instant (i.e., if the last component is the day, it represents a whole day; if it's the hour, it represents the whole hour.) So `2014/1/31` represents one whole day. `2014/1/31-2014/2/1` represents two days. `2014/1/31T14:50-2014/1/31T14:51` represents two whole minutes. `2014/1/31T12-` represents all time from 2014/1/31 noon forward. The options in Table C.7 are specified with time ranges.

**Country codes:** The `--scc`, `--dcc`, and `--any-cc` parameters take a comma-separated list of two-letter country codes, as specified by IANA.<sup>20</sup> There are also four special codes: `--` for unknown, `a1` for anonymous proxy, `a2` for satellite provider, and `o1` for other. The options in Table C.8 are specified with country codes. These options require a country-code mapping file to be built and installed.

**Address types:** The `--stype` and `--dtype` parameters accept a single index. Records pass if their source or destination addresses produce a matching index when looked up in address-types mapping file. These parameters require the mapping file to be built and installed.

## rwfilter Output Parameters

Table C.10: `rwfilter` Output Parameters

Parameter	Example	Description
<code>--pass-destination</code>	<code>stdout</code>	Send SiLK flow records matching all partitioning parameters to pipe or file
<code>--fail-destination</code>	<code>faildata.rw</code>	Like <code>--pass</code> , but for records failing to match
<code>--all-destination</code>	<code>allrecs.rw</code>	Like <code>--pass</code> , but for all records
<code>--print-statistics</code>	—	Print (default to <code>stderr</code> ) count of records passing and failing
<code>--print-volume-statistics</code>	<code>out-vol.txt</code>	Print counts of flows/bytes/packets read, passing, and failing to named file or <code>stderr</code>

<sup>20</sup><https://www.iana.org/domains/root/db/>

### Miscellaneous *rwfilter* Parameters

Additional *rwfilter* parameters that are commonly useful in analysis or maintaining the repository are listed below. These depend on the implementation and are described in Table C.11.

The `--threads` parameter takes an integer scalar  $N$  to specify using  $N$  threads to read input files and filter records. The default value is one or the value of the `SILK_RWFILTER_THREADS` environment variable if that is set. Using multiple threads is preferable for queries that look at many files but return few records. Current experience is that performance peaks at about two to three threads per CPU (core) on the host running *rwfilter*, but this result is variable with the type of query and the number of records returned from each file. There is no advantage in having more threads than input files. There may also be a point of diminishing returns, which seems to be around 20 threads.

To improve query efficiency when few records are needed, the `--max-pass-records` parameter allows the analyst to specify the maximum number of records to return via the path specified by the `--pass` parameter.

Table C.11: Miscellaneous *rwfilter* Parameters

Parameter	Description
<code>--print-missing-files</code>	Print names of missing repository files to stderr. Doubles as a selection parameter.
<code>--threads</code>	Specify number of process threads to be used in filtering
<code>--max-pass-records</code>	Specifies the maximum number of records to return as matching partitioning parameters
<code>--max-fail-records</code>	Specifies the maximum number of records to return as <i>not</i> matching partitioning parameters
<code>--dry-run</code>	Performs a sanity check on parameters. Does not retrieve records. Prints a list of the names of files that would be accessed.

For additional parameters, see Table C.17 and Table C.18 .

## C.4 `rwstats` Command Summary

### `rwstats`

<b>Description</b>	Summarizes SiLK flow records by one of a limited number of key-and-value pairs and displays the results as a top- <i>N</i> or bottom- <i>N</i> list
<b>Call</b>	<code>rwstats flowrecs.rw --fields=protocol --values=Records --top --count=20</code>
<b>Parameters</b>	<p>Choose one or none:</p> <ul style="list-style-type: none"> <li><code>--top</code> Prints the top <i>N</i> keys and their values (default)</li> <li><code>--bottom</code> Prints the bottom <i>N</i> keys and their values</li> <li><code>--overall-stats</code> Prints minima, maxima, quartiles, and interval-count statistics for bytes, packets, and bytes per packet across all flows</li> <li><code>--detail-proto-stats</code> Prints overall statistics for each specified protocol. Protocols are specified as integers or ranges separated by commas.</li> </ul> <p>Choose one for <code>--top</code> or <code>--bottom</code>:</p> <ul style="list-style-type: none"> <li><code>--count</code> Displays the specified number of key-and-value pairs</li> <li><code>--percentage</code> Displays key-and-value pairs where the value is greater than (<code>--top</code>) or less than (<code>--bottom</code>) this percentage of the total value</li> <li><code>--threshold</code> Displays key-and-value pairs where the the specified constant is greater than or less than a threshold value.</li> </ul> <p>Options for <code>--top</code> or <code>--bottom</code>:</p> <ul style="list-style-type: none"> <li><code>--fields</code> Uses the indicated fields as the key (see Table C.13) – required</li> <li><code>--values</code> Calculates values for the specified fields (default: Records)</li> <li><code>--presorted-input</code> Assumes input is already sorted by key</li> <li><code>--no-percents</code> Does not display percent-of-total or cumulative-percentage</li> <li><code>--bin-time</code> Adjusts <code>sTime</code> and <code>eTime</code> to multiple of the argument in seconds</li> <li><code>--temp-directory</code> Specifies the location for temporary data when memory is exceeded</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.5 *rwcount* Command Summary

### **rwcount**

<b>Description</b>	Calculates volumes over time periods of equal duration
<b>Call</b>	<code>rwcount flowrecs.rw --bin-size=3600</code>
<b>Parameters</b>	<p><code>--bin-size</code> Specifies the number of seconds per bin (default 30)</p> <p><code>--load-scheme</code> Specifies how the flow volume is allocated to bins (see Table C.12 for details)</p> <p><code>--skip-zeroes</code> Does not print empty bins</p> <p><code>--start-time</code> Specifies the initial time of the first bin</p> <p><code>--end-time</code> Specifies a time in the last bin, extended to make a whole bin</p> <p><code>--bin-slots</code> Displays timestamps as internal bin indices</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

Table C.12: Time distribution options for *rwcount --load-scheme*

Value	Parameter Name	Volume Allocation	Guidelines for Use
0	<code>bin-uniform</code>	Allocates equal parts of volume to every bin in timespan	Computes the average load per bin, smoothing out peaks and valleys.
1	<code>start-spike</code>	Stores entire volume in the first millisecond of flow (i.e., the flow's <code>stime</code> bin)	Emphasizes the onset of periodic behavior. Puts all packets and bytes into one bin even if the flow spans multiple bins.
2	<code>end-spike</code>	Stores entire volume in the last millisecond of the flow (i.e., the flow's <code>etime</code> bin)	Emphasizes flow termination. Puts all packets and bytes into one bin even if the flow spans multiple bins.
3	<code>middle-spike</code>	Stores entire volume in the middle millisecond of the flow	Emphasizes payload transfer. Puts all packets and bytes into one bin even if the flow spans multiple bins.
4	<code>time-proportional</code>	Proportionally allocates the flow's bytes, packets, and record count (1 for one flow) to all bins in the flow according to how much time the flow spent in each bin's timespan	<b>Default load scheme; recommended for most analyses.</b> Gives the average load per time period. Smooths out peaks and valleys over time.
5	<code>max-volume</code>	Assigns entire flow volume to each bin	Overestimates load; computes worst-case scenario for service loading.
6	<code>min-volume</code>	Assigns one flow to each bin	Underestimates load; computes best case scenario for service loading.

- For `--load-scheme` values 0 through 4, the flow record count adds up to 1. The byte and packet counts add up to the counts in the flow record.
- For `--load-scheme` values 5 and 6, the flow record count does not add up to 1. The byte and packet counts do not add up to the counts in the flow record.

## C.6 `rwcut` Command Summary

### `rwcut`

<b>Description</b>	Reads SiLK flow data and displays it as text
<b>Call</b>	<code>rwcut flowrecs.rw --fields=1-5,sTime</code>
<b>Parameters</b>	<p>Choose one or none:</p> <ul style="list-style-type: none"> <li><code>--fields</code> Specifies which fields to display (default is 1-12)</li> <li><code>--all-fields</code> Displays all fields</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><code>--start-rec-num</code> Specifies record offset of first record from start of file</li> <li><code>--end-rec-num</code> Specifies record offset of last record from start of file</li> <li><code>--tail-recs</code> Specifies record offset of first record from end of file (cannot combine with <code>--start-rec-num</code> or <code>--end-rec-num</code>)</li> <li><code>--num-recs</code> Specifies maximum number of output records</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>



Table C.13: Arguments for the `--fields` Parameter

Field Number	Field Name	Description
1	sIP	Source IP address for flow
2	dIP	Destination IP address for flow
3	sPort	Source port for flow (or 0)
4	dPort	Destination port for flow (or 0)
5	protocol	Transport-layer protocol number for flow
6	packets, pkts	Number of packets in flow
7	bytes	Number of bytes in flow (starting with IP header)
8	flags	Cumulative TCP flag fields of flow (or blank)
9	sTime	Start date and time of flow
10	duration	Duration of flow
11	eTime	End date and time of flow
12	sensor	Sensor that collected the flow
13	in	Ingress interface or VLAN on sensor (usually zero)
14	out	Egress interface or VLAN on sensor (usually zero)
15	nhIP	Next-hop IP address (usually zero)
16	sType	Type of source IP address (pmap required)
17	dType	Type of destination IP address (pmap required)
18	scc	Source country code (pmap required)
19	dcc	Destination country code (pmap required)
20	class	Class of sensor that collected flow
21	type	Type of flow for this sensor class
—	iType	ICMP type for ICMP and ICMPv6 flows (SiLK V3.8.1+)
—	iCode	ICMP code for ICMP and ICMPv6 flows (SiLK V3.8.1+)
25	icmpTypeCode	Both ICMP type and code values (before SiLK V3.8.1)
26	initialFlags	TCP flags in initial packet
27	sessionFlags	TCP flags in remaining packets
28	attributes	Termination conditions
29	application	Standard port for application that produced the flow

## C.7 rwsort Command Summary

### rwsort

<b>Description</b>	Sorts SiLK flow records using key field(s)
<b>Call</b>	rwsort unsorted1.rw unsorted2.rw --fields=1,3 --output-path=sorted.rw
<b>Parameters</b>	<p><b>--fields</b> Specifies key fields for sorting (required)</p> <p><b>--presorted-input</b> Specifies only merging of already sorted input files</p> <p><b>--reverse</b> Specifies sort in descending order</p> <p><b>--temp-directory</b> Specifies location of high-speed storage for temp files</p> <p><b>--sort-buffer-size</b> Specifies the in-memory sort buffer (2 GB default)</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.8 `rwuniq` Command Summary

### `rwuniq`

<b>Description</b>	Counts records per combination of multiple-field keys
<b>Call</b>	<code>rwuniq filterfile.rw --fields=1-5,sensor --values=Records</code>
<b>Parameters</b>	<p><code>--fields</code> Specifies fields to use as key (required)</p> <p><code>--values</code> Specifies summary counts (default: Records)</p> <p><code>--bin-time</code> Establishes bin size for time-oriented bins</p> <p><code>--presorted-input</code> Reduces memory requirements for presorted records</p> <p><code>--sort-output</code> Sorts results by key, as specified in the <code>--fields</code> parameter</p> <p>For options to filter output rows, see Table C.14. For additional parameters, see Table C.17 and Table C.18.</p>

Table C.14: Output-Filtering Options for `rwuniq`

Parameter	Description
<code>--bytes</code>	Only output rows whose byte counts are in the specified range
<code>--packets</code>	Only output rows total packet counts are in the specified range
<code>--flows</code>	Only output rows whose flow (record) counts are in the specified range
<code>--sip-distinct</code>	Only output rows whose counts of distinct (unique) source IP addresses are in the specified range
<code>--dip-distinct</code>	Only output rows whose counts of distinct (unique) destination IP addresses are in the specified range

## C.9 rwnetmask Command Summary

### rwnetmask

<b>Description</b>	Zeroes all bits past the specified prefix length on the specified address in SiLK flow records
<b>Call</b>	<code>rwnetmask flows.rw --4sip-prefix-length=24 --output-path=sip-24.rw</code>
<b>Parameters</b>	<p>Choose one or more:</p> <ul style="list-style-type: none"> <li><code>--4sip-prefix-length</code> Gives number of high bits of source IPv4 address to keep</li> <li><code>--4dip-prefix-length</code> Gives number of high bits of destination IPv4 to keep</li> <li><code>--4nhip-prefix-length</code> Gives number of high bits of next-hop IPv4 to keep</li> <li><code>--6sip-prefix-length</code> Gives number of high bits of source IPv6 to keep</li> <li><code>--6dip-prefix-length</code> Gives number of high bits of destination IPv6 to keep</li> <li><code>--6nhip-prefix-length</code> Gives number of high bits of next-hop IPv6 to keep</li> </ul> <p>For additional parameters, see Table <a href="#">C.17</a> and Table <a href="#">C.18</a>.</p>

## C.10 rwcatt Command Summary

### rwcatt

<b>Description</b>	Concatenates SiLK flow record files and copies to a new file
<b>Call</b>	<code>rwcatt someflows.rw moreflows.rw --output-path=allflows.rw</code>
<b>Parameters</b>	<code>--ipv4-output</code> Converts IPv6 records to IPv4 records following the <code>asv4</code> policy; ignores other IPv6 records <code>--byte-order</code> Writes the output in this byte order. Possible choices are <code>native</code> (the default), <code>little</code> , and <code>big</code>

For additional parameters, see Table [C.17](#) and Table [C.18](#).

## C.11 `rwappend` Command Summary

### `rwappend`

<b>Description</b>	Appends the flow records from the successive files to the first file
<b>Call</b>	<code>rwappend allflows.rw laterflows.rw</code>
<b>Parameters</b>	<p><code>--create</code> Creates the output file if it does not already exist. Determines the format and version of the output file from the flow record file optionally named in this parameter</p> <p><code>--print-statistics</code> Prints to standard error the count of records that are read from each input file and written to the output file</p> <p>For additional parameters, see Table <a href="#">C.17</a> and Table <a href="#">C.18</a>.</p>

## C.12 rwsplit Command Summary

### rwsplit

<b>Description</b>	Divides the flow records into successive files
<b>Call</b>	<code>rwsplit allflows.rw --flow-limit=1000 --basename=sample</code>
<b>Parameters</b>	<p>Choose one:</p> <ul style="list-style-type: none"> <li><code>--ip-limit</code> Specifies the IP address count at which to begin a new sample file</li> <li><code>--flow-limit</code> Specifies the flow count at which to begin a new sample file</li> <li><code>--packet-limit</code> Specifies the packet count at which to begin a new sample file</li> <li><code>--byte-limit</code> Specifies the byte count at which to begin a new sample file</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><code>--basename</code> Specifies the base name for output sample files (required)</li> <li><code>--sample-ratio</code> Specifies the denominator for the ratio of records read to number written in a sample file (e.g., 100 means to write 1 out of 100 records)</li> <li><code>--seed</code> Seeds the random number generator with this value</li> <li><code>--file-ratio</code> Specifies the denominator for the ratio of sample filenames generated to the total number written (e.g., 10 means 1 of every 10 files will be saved)</li> <li><code>--max-outputs</code> Specifies the maximum number of files to write to disk</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.13 rwtuc Command Summary

### rwtuc

<b>Description</b>	Generates SiLK flow records from textual input similar to rwcut output
<b>Call</b>	<code>rwtuc flows.rw.txt --output-path=flows.rw</code>
<b>Parameters</b>	<p><b>--fields</b> Specifies the fields to parse from the input. Values may be</p> <ul style="list-style-type: none"> <li>• a field number: 1–15, 20–21, 26–29</li> <li>• a field name equivalent to one of the field numbers above (see Table C.13 on page 193)</li> <li>• the keyword <code>ignore</code> for an input column not to be included in the output records.</li> </ul> <p>The parameter is unnecessary if the input file has appropriate column headings.</p> <p><b>--bad-input-lines</b> Specifies the file or stream to write each bad input line to (filename and line number prepended)</p> <p><b>--verbose</b> Prints an error message for each bad input line to standard error</p> <p><b>--stop-on-error</b> Prints an error message for a bad input line to standard error and exits</p> <p><b>--fixed-value-parameter</b> Uses the value as a fixed value for this field in all records. See Table C.15 for the parameter name for each field.</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

Table C.15: Fixed-Value Parameters for rwtuc

Field	Parameter Name	Field	Parameter Name
sIP	<code>--saddress</code>	in	<code>--input-index</code>
dIP	<code>--daddress</code>	out	<code>--output-index</code>
sPort	<code>--sport</code>	nhIP	<code>--next-hop-ip</code>
dPort	<code>--dport</code>	class	<code>--class</code>
protocol	<code>--protocol</code>	type	<code>--type</code>
packets	<code>--packets</code>	iType	<code>--icmp-type</code>
bytes	<code>--bytes</code>	iCode	<code>--icmp-code</code>
flags	<code>--flags-all</code>	initialFlags	<code>--flags-initial</code>
sTime	<code>--stime</code>	sessionFlags	<code>--flags-session</code>
duration	<code>--duration</code>	attributes	<code>--attributes</code>
eTime	<code>--etime</code>	application	<code>--application</code>
sensor	<code>--sensor</code>		



## C.14 `rwset` Command Summary

### `rwset`

<b>Description</b>	Generates IP-set files from flows
<b>Call</b>	<code>rwset flows.rw --sip-file=flows_sip.set</code>
<b>Parameters</b>	<p>Choose one or more:</p> <ul style="list-style-type: none"> <li><code>--sip-file</code> Specifies an IP-set file to generate with source IP addresses from the flows records</li> <li><code>--dip-file</code> Specifies an IP-set file to generate with destination IP addresses from the flows records</li> <li><code>--nhip-file</code> Specifies an IP-set file to generate with next-hop IP addresses from the flows records</li> <li><code>--any-file</code> Specifies an IP-set file to generate with both source and destination IP addresses from the flows records</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><code>--record-version</code> Specifies the version of records to write to a file. Allowed arguments are 0, 2, 3, and 4; record version affects size and backwards compatibility. It can also be set by the <code>SILK_IPSET_RECORD_VERSION</code> environment variable.</li> <li><code>--invocation-strip</code> Does not copy command lines from the input files to the output files</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.15 `rwsetcat` Command Summary

### `rwsetcat`

<b>Description</b>	Lists contents of IP-set files as text on standard output
<b>Call</b>	<code>rwsetcat low_sip.set &gt;low_sip.set.txt</code>
<b>Parameters</b>	<p><b>--network-structure</b> Prints the network structure of the set.            Syntax: <code>[v6: v4:][list-lengths[S][/[summary-lengths]]]</code>            A length may be expressed as an integer prefix length (often preferred) or a letter: T for total address space (/0), A for /8, B for /16, C for /24, X for /27, and H for Host (/32 for IPv4, /128 for IPv6); with S for default summaries.</p> <p><b>--cidr-blocks</b> Groups IP addresses that fill a CIDR block</p> <p><b>--ip-ranges</b> Groups consecutive addresses; provides the most compact display</p> <p><b>--count-ips</b> Prints the number of IP addresses. Disables default printing of addresses</p> <p><b>--print-statistics</b> Prints set statistics (min-/max-IP address, etc.). Also disables default printing of addresses</p> <p><b>--print-ips</b> Prints IP addresses when count or statistics parameter is given</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.16 `rwsettool` Command Summary

### `rwsettool`

<b>Description</b>	Manipulates IP-set files to produce new IP-set files
<b>Call</b>	<code>rwsettool --intersect src1.set src2.set --output-path=both.set</code>
<b>Parameters</b>	<p>Choose one:</p> <ul style="list-style-type: none"> <li><b>--union</b> Creates set containing IP addresses in any input file</li> <li><b>--intersect</b> Creates set containing IP addresses in all input files</li> <li><b>--difference</b> Creates set containing IP addresses from first file not in any of the remaining files</li> <li><b>--mask</b> Creates set containing one IP address from each block of the specified bitmask length when any of the input IPsets have an IP address in that block</li> <li><b>--fill-blocks</b> Creates an IPset containing a completely full block with the specified prefix length when any of the input IPsets have an IP address in that block</li> <li><b>--sample</b> Creates an IPset containing a random sample of IP addresses from all input IPsets. Requires either the <b>--size</b> or <b>--ratio</b> option</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><b>--size</b> Specifies the sample size (number of IP addresses sampled from each input IPset). Requires the <b>--sample</b> parameter</li> <li><b>--ratio</b> Specifies the probability (as a floating point value between 0.0 and 1.0) that an IP address will be sampled. Requires the <b>--sample</b> parameter</li> <li><b>--seed</b> Specifies the random number seed integer for the <b>--sample</b> parameter and is used only with that parameter</li> <li><b>--note-strip</b> Does not copy notes from the input files to the output file</li> <li><b>--invocation-strip</b> Does not copy command history from the input files to the output file</li> <li><b>--record-version</b> Specifies the IP-set record version to write. v2 only supports IPv4; v3 requires SiLK 3; v4 is more compact and requires SiLK 3.7 or later; 0 uses v2 for IPv4 sets, and v3 otherwise.</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.17 `rwsetbuild` Command Summary

### `rwsetbuild`

<b>Description</b>	Create a binary IPset file from a list of IP addresses
<b>Call</b>	<code>rwsetbuild myset.set.txt myset.set</code>
<b>Parameters</b>	<code>--ip-ranges=DELIM</code> Range of IP addresses, where DELIM is the delimiter. <code>--invocation-strip</code> Do not include the command line in the file
For additional parameters, see Table <a href="#">C.17</a> and Table <a href="#">C.18</a> .	

## C.18 rwbag Command Summary

### rwbag

<b>Description</b>	Generates bags from flow records		
<b>Call</b>	<code>rwbag flow.rw --sip-bytes=x.bag --sip-flows=y.bag</code>		
<b>Parameters</b>	<b>--key-count</b> Writes bag of <i>count</i> by unique <i>key</i> value. May be specified multiple times. Allowed values for <i>key</i> and <i>count</i> are		
	<b>Key</b>		<b>Count</b>
	sip	Source IP address	flows
	dip	Destination IP address	packets
	nhip	Next-hop IP address	bytes
	input	Router input port	
	output	Router output port	
	sport	Source port	
	dport	Destination port	
	proto	Protocol	
	sensor	Sensor ID	
	For additional parameters, see Table C.17 and Table C.18.		

## C.19 `rwbagbuild` Command Summary

### `rwbagbuild`

<b>Description</b>	Creates a binary bag from non-flow data (expressed as text)
<b>Call</b>	<code>rwbagbuild --bag-input=ip-byte.bag.txt --key-type=sIPv4 --counter-type=sum-bytes --output-path=ip-byte.bag</code>
<b>Parameters</b>	<p>Choose one:</p> <ul style="list-style-type: none"> <li><b>--set-input</b> Creates a bag from the specified IPset, which may be <code>stdin</code> or a hyphen (-)</li> <li><b>--bag-input</b> Creates a bag from a delimiter-separated text file, which can be <code>stdin</code> or a hyphen (-)</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><b>--delimiter</b> Specifies the delimiter separating the key and value for the <code>--bag-input</code> parameter. Cannot be the pound sign (#) or the line-break character (new line)</li> <li><b>--default-count</b> Specifies the integer count for each key in the new bag, overriding any values present in the input</li> <li><b>--key-type</b> Sets the key type to this value. Allowable options are shown in Table C.16.</li> <li><b>--counter-type</b> Sets the counter type to this value. Allowable options are shown in Table C.16.</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

Table C.16: rwbagbuild Key or Value Options

Type	Description, allowed values
sIPv4	Source IP addresses, IPv4 only, dotted quad, CIDR, wildcard, or integer
dIPv4	Destination IP address, IPv4 only, dotted quad, CIDR, wildcard, or integer
sPort	Source port, integer 0–65,535
dPort	Destination port, integer 0–65,535
protocol	IP protocol, integer 0–255
packets	Packet count, integer
bytes	Byte count, integer
flags	Bit string of TCP cumulative TCP flags (CEUAPRSF), integer 0–255
sTime	Starting time of the flow, integer seconds from UNIX epoch
duration	Duration of the flow, integer seconds
eTime	Ending time of the flow, integer seconds from UNIX epoch
sensor	Sensor ID, integer
input	SNMP index of input interface, integer
output	SNMP index of output interface, integer
nhIPv4	Next-hop IP address, IPv4 only, dotted quad, CIDR, wildcard, or integer
initialFlags	TCP flags in first packet in the flow, integer 0–255
sessionFlags	Cumulative TCP flags excluding the first packet, integer 0–255
attributes	Flags for termination conditions and packet size uniformity, integer
application	Guess as to the content of the flow, as set by the flow generator, integer 0–65,535
class	Class of the sensor, integer
type	Type of the flow, integer
icmpTypeCode	An encoded version of the ICMP type and code, where the type is in the upper byte and the code is in the lower byte
sIPv6	Source IP address, IPv6, canonical form or integer
dIPv6	Destination IP address, IPv6, canonical form or integer
nhIPv6	Next-hop IP address, IPv6, canonical form or integer
records	Count of flows, integer
sum-packets	Sum of packet counts, integer
sum-bytes	Sum of byte counts, integer
sum-duration	Sum of duration values, integer seconds
any-IPv4	Source, destination, or next-hop IPv4 address, dotted quad or integer
any-IPv6	Source, destination or next-hop IPv6 address, canonical form or integer
any-port	Source or destination port, integer 0–65,535
any-snmp	Input or output SNMP index of interface, integer
any-time	Start or end time value, integer seconds since UNIX epoch
custom	An integer

## C.20 `rwbagcat` Command Summary

### `rwbagcat`

<b>Description</b>	Displays or summarizes bag contents
<b>Call</b>	<code>rwbagcat x.bag --output-path=x.bag.txt</code>
<b>Parameters</b>	<p>Choose one or none:</p> <ul style="list-style-type: none"> <li><b>--network-structure</b> Prints the sum of counters for each specified CIDR block in the comma-separated list of CIDR block sizes and/or letters</li> <li><b>--bin-ips</b> Inverts the bag and counts keys by distinct volume values. Allowed arguments are <b>linear</b> (count keys with each volume [the default]), <b>binary</b> (count keys with volumes that fall in ranges based on powers of 2), and <b>decimal</b> (count keys that fall in ranges determined by a decimal logarithmic scale).</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><b>--mincounter</b> Displays only entries with counts of at least the value given as the argument</li> <li><b>--maxcounter</b> Displays only entries with counts no larger than the value given as the argument</li> <li><b>--minkey</b> Displays only entries with keys of at least the value given as the argument</li> <li><b>--maxkey</b> Displays only entries with keys no larger than the value given as the argument</li> <li><b>--key-format</b> Specifies the formatting of keys for display. Allowed arguments are <b>canonical</b> (display keys as IP addresses in canonical format), <b>zero-padded</b> (display keys as IP addresses with zeroes added to fully fill width of column), <b>decimal</b> (display keys as decimal integer values), <b>hexadecimal</b> (display keys as hexadecimal integer values), and <b>force-ipv6</b> (display all keys as IP addresses in the canonical form for IPv6 with no IPv4 notation).</li> <li><b>--mask-set</b> Outputs records whose keys appear in the argument set-file</li> <li><b>--zero-counts</b> Prints keys with a counter of zero (requires <b>--mask-set</b> or both <b>--minkey</b> and <b>--maxkey</b>)</li> <li><b>--print-statistics</b> Prints statistics about the bag to given file or standard output</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>



C.21 `rwbagtool` Command Summary**rwbagtool**

<b>Description</b>	Manipulates bags and generates cover sets
<b>Call</b>	<code>rwbagtool --add x.bag y.bag --output-path=z.bag</code>
<b>Parameters</b>	<p>Choose one (or none with one input bag):</p> <ul style="list-style-type: none"> <li><code>--add</code> Adds bags together (union)</li> <li><code>--subtract</code> Subtracts from the first bag all the other bags (difference)</li> <li><code>--minimize</code> Writes to the output the minimum counter for each key across all input bags</li> <li><code>--maximize</code> Writes to the output the maximum counter for each key across all input bags</li> <li><code>--divide</code> Divides the first bag by the second bag</li> <li><code>--scalar-multiply</code> Multiplies each counter in the bag by the specified value. Accepts a single bag file as input.</li> <li><code>--compare</code> Compares key/value pairs in exactly two bag files using the operation (OP) specified by the argument. Keeps only those keys in the first bag that also appear in the second bag and whose counter satisfies the OP relation with those in the second bag. The counter for each key that remains is set to 1. Allowed OPs are <code>lt</code> (less than), <code>le</code> (less than or equal to), <code>eq</code> (equal to), <code>ge</code> (greater than or equal to), <code>gt</code> (greater than).</li> </ul> <p>Choose zero or more masking/limiting parameters to restrict the results of the above operation or the sole input bag:</p> <ul style="list-style-type: none"> <li><code>--intersect</code> Intersects the specified set with keys in the bag</li> <li><code>--complement-intersect</code> Masks keys in the bag using IP addresses <i>not</i> in given IP-set file</li> <li><code>--minkey</code> Cuts bag to entries with key of at least the value given as an argument</li> <li><code>--maxkey</code> Cuts bag to entries with key of at most the value given as an argument</li> <li><code>--mincounter</code> Outputs records whose counter is at least the value given as an argument, an integer</li> <li><code>--maxcounter</code> Outputs records whose counter is not more than the value given as an argument, an integer</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><code>--coverset</code> Generates an IPset for bag keys instead of creating a bag</li> <li><code>--invert</code> Counts keys for each unique counter value</li> <li><code>--note-strip</code> Does not copy notes from the input files to the output file</li> <li><code>--output-path</code> Specifies where resulting bag or set should be stored</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.22 `rwfileinfo` Command Summary

### `rwfileinfo`

<b>Description</b>	Displays summary information about one or more SiLK files																																																							
<b>Call</b>	<code>rwfileinfo allflows.rw --fields=count-records,command-lines</code>																																																							
<b>Parameters</b>	<p><code>--fields</code> Selects which summary information to display via number or name (by default, all the available fields). Possible values include</p> <table border="1"> <thead> <tr> <th>#</th> <th>Field</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td>1</td> <td><code>format</code></td> <td>Binary file format indicator</td> </tr> <tr> <td>2</td> <td><code>version</code></td> <td>Version of file header</td> </tr> <tr> <td>3</td> <td><code>byte-order</code></td> <td>Byte order of words written to disk</td> </tr> <tr> <td>4</td> <td><code>compression</code></td> <td>Type of space compression used</td> </tr> <tr> <td>5</td> <td><code>header-length</code></td> <td>Number of bytes in file header</td> </tr> <tr> <td>6</td> <td><code>record-length</code></td> <td>Number of bytes in fixed-length records</td> </tr> <tr> <td>7</td> <td><code>count-records</code></td> <td>Number of records in the file unless <code>record-length=1</code></td> </tr> <tr> <td>8</td> <td><code>file-size</code></td> <td>Total number of bytes in the file on disk</td> </tr> <tr> <td>9</td> <td><code>command-lines</code></td> <td>List of stored commands that generated this file</td> </tr> <tr> <td>10</td> <td><code>record-version</code></td> <td>Version of records in file</td> </tr> <tr> <td>11</td> <td><code>silk-version</code></td> <td>Software version of SiLK tool that produced this file</td> </tr> <tr> <td>12</td> <td><code>packed-file-info</code></td> <td>Information from packing process</td> </tr> <tr> <td>13</td> <td><code>probe-name</code></td> <td>Probe info for files created by flowcap</td> </tr> <tr> <td>14</td> <td><code>annotations</code></td> <td>List of notes</td> </tr> <tr> <td>15</td> <td><code>prefix-map</code></td> <td>Prefix map name and header version</td> </tr> <tr> <td>16</td> <td><code>ipset</code></td> <td>IP-set format information</td> </tr> <tr> <td>17</td> <td><code>bag</code></td> <td>Bag key and count information</td> </tr> </tbody> </table> <p><code>--summary</code> Prints a summary of total files, file sizes, and records</p>		#	Field	Description	1	<code>format</code>	Binary file format indicator	2	<code>version</code>	Version of file header	3	<code>byte-order</code>	Byte order of words written to disk	4	<code>compression</code>	Type of space compression used	5	<code>header-length</code>	Number of bytes in file header	6	<code>record-length</code>	Number of bytes in fixed-length records	7	<code>count-records</code>	Number of records in the file unless <code>record-length=1</code>	8	<code>file-size</code>	Total number of bytes in the file on disk	9	<code>command-lines</code>	List of stored commands that generated this file	10	<code>record-version</code>	Version of records in file	11	<code>silk-version</code>	Software version of SiLK tool that produced this file	12	<code>packed-file-info</code>	Information from packing process	13	<code>probe-name</code>	Probe info for files created by flowcap	14	<code>annotations</code>	List of notes	15	<code>prefix-map</code>	Prefix map name and header version	16	<code>ipset</code>	IP-set format information	17	<code>bag</code>	Bag key and count information
#	Field	Description																																																						
1	<code>format</code>	Binary file format indicator																																																						
2	<code>version</code>	Version of file header																																																						
3	<code>byte-order</code>	Byte order of words written to disk																																																						
4	<code>compression</code>	Type of space compression used																																																						
5	<code>header-length</code>	Number of bytes in file header																																																						
6	<code>record-length</code>	Number of bytes in fixed-length records																																																						
7	<code>count-records</code>	Number of records in the file unless <code>record-length=1</code>																																																						
8	<code>file-size</code>	Total number of bytes in the file on disk																																																						
9	<code>command-lines</code>	List of stored commands that generated this file																																																						
10	<code>record-version</code>	Version of records in file																																																						
11	<code>silk-version</code>	Software version of SiLK tool that produced this file																																																						
12	<code>packed-file-info</code>	Information from packing process																																																						
13	<code>probe-name</code>	Probe info for files created by flowcap																																																						
14	<code>annotations</code>	List of notes																																																						
15	<code>prefix-map</code>	Prefix map name and header version																																																						
16	<code>ipset</code>	IP-set format information																																																						
17	<code>bag</code>	Bag key and count information																																																						
	For additional parameters, see Table C.17 and Table C.18.																																																							

## C.23 `rwpmapbuild` Command Summary

### `rwpmapbuild`

<b>Description</b>	Creates a prefix map from a text file
<b>Call</b>	<code>rwpmapbuild --input-file=sample.pmap.txt --output-file=sample.pmap</code>
<b>Parameters</b>	<p><b>--input-file</b> Specifies the text file that contains the mapping between addresses and prefixes. When omitted, read from <code>stdin</code></p> <p><b>--mode</b> Specifies the type of input as if a <code>mode</code> statement appeared in the input. Valid values are <code>ipv4</code>, <code>ipv6</code>, and <code>proto-port</code>.</p> <p><b>--output-file</b> Specifies the filename for the binary prefix map file</p> <p><b>--ignore-errors</b> Writes the output file despite any errors in the input</p> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.24 `rwpmlookup` Command Summary

### `rwpmlookup`

<b>Description</b>	Shows label associated with an addresses' prefix map, address type, or country code
<b>Call</b>	<code>rwpmlookup --map-file=spyware.map address-list.txt</code>
<b>Parameters</b>	<p>Choose one:</p> <ul style="list-style-type: none"> <li><b>--map-file</b> Specifies the pmap that contains the prefix map to query</li> <li><b>--address-types</b> Finds IP addresses in the address-types mapping file specified in the argument or in the default file when no argument is provided</li> <li><b>--country-codes</b> Finds IP addresses in the country-code mapping file specified in the argument or in the default file when no argument is provided</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><b>--fields</b> Specifies the fields to print. Allowed values are: <b>key</b> (key used to query), <b>value</b> (label from prefix map), <b>input</b> (the text read from the input file [excluding comments] or IP address in canonical form from set input file), <b>block</b> (CIDR block containing the key), <b>start-block</b> (low address in CIDR block containing the key), and <b>end-block</b> (high address in CIDR block containing the key.) Default is <b>key,value</b>.</li> <li><b>--no-files</b> Does not read from files and instead treat the command line arguments as the IP addresses or protocol/port pairs to find</li> <li><b>--no-errors</b> Does not report errors parsing the input</li> <li><b>--ipset-files</b> Treats the command-line arguments as names of binary IP-set files to read</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.25 rwmatch Command Summary

### rwmatch

<b>Description</b>	Matches IPv4 flow records that have stimulus-response relationships (IPv6 support introduced in Version 3.9.0)
<b>Call</b>	<code>rwmatch --relate=1,2 --relate=2,1 query.rw response.rw matched.rw</code>
<b>Parameters</b>	<p>Choose one or none:</p> <ul style="list-style-type: none"> <li><b>--absolute-delta</b> Includes potentially matching flows that start less than the interval specified by <b>--time-delta</b> after the end of the initial flow of the current match (default)</li> <li><b>--relative-delta</b> Continues matching with flows that start within the interval specified by <b>--time-delta</b> from the greatest end time seen for previous members of the current match</li> <li><b>--infinite-delta</b> After forming the initial pair of the match, continues matching on relate fields alone, ignoring time</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><b>--relate</b> Specifies the numeric field IDs (1–8; see Figure C.13 on page 193) that identify stimulus and response (required; may be specified multiple times) Starting with Version 3.9.0 values may be 1–8, 12–14, 20–21, 26–29, iType, and iCode; values may be specified by name or numeric ID.</li> <li><b>--time-delta</b> Specifies the number of seconds by which a time window is extended beyond a record's end time. The default value is zero.</li> <li><b>--symmetric-delta</b> Also makes an initial match for a query that starts between a response's start time and its end time extended by <b>--time-delta</b></li> <li><b>--unmatched</b> Includes unmatched records from the query file and/or the response file in the output. Allowed arguments (case-insensitive) are one of these: <b>q</b> (query file), <b>r</b> (response file), <b>b</b> (both)</li> </ul> <p>For additional parameters, see Table C.17 and Table C.18.</p>

## C.26 `rwgroup` Command Summary

### `rwgroup`

<b>Description</b>	Flags flow records that have common attributes
<b>Call</b>	<code>rwgroup --id-fields=sIP --delta-field=sTime --delta-value=3600 --output-path=grouped.rw</code>
<b>Parameters</b>	<p>Choose one or both:</p> <ul style="list-style-type: none"> <li><code>--id-fields</code> Specifies the fields that need to be identical</li> <li><code>--delta-field</code> Specifies the field that needs to be close. Requires the <code>--delta-value</code> parameter</li> </ul> <p>Options:</p> <ul style="list-style-type: none"> <li><code>--delta-value</code> Specifies closeness for the <code>--delta-field</code> parameter</li> <li><code>--objective</code> Specifies that the <code>--delta-value</code> argument applies relative to the first record, rather than the most recent</li> <li><code>--rec-threshold</code> Specifies the minimum number of records in a group</li> <li><code>--summarize</code> Produces a single record as output for each group, rather than all flow records</li> </ul> <p>For additional parameters, see Table <a href="#">C.17</a> and Table <a href="#">C.18</a>.</p>

## C.27 Features Common to Several Commands

Many of the SiLK tools share features such as using common parameters, providing help, handling the two versions of IP addresses, and controlling the overwriting of existing output files.

### C.27.1 Parameters Common to Several Commands

Many options apply to several of the SiLK tools, as shown in Table C.17.

Table C.18 lists the same options as in Table C.17 and provides descriptions of the options. Three of the options described accept a small number of fixed values; acceptable values for `--ip-format` are listed and described in Table C.19, values for `--timestamp-format` are described in Table C.20, and values for `--ipv6-policy` are described in Table C.18 on page 217.

Table C.17: Common Parameters in Essential SiLK Tools

Parameter	rwfilter	rwstats	rwcount	rwcut	rwsort	rwuniq
--help	✓	✓	✓	✓	✓	✓
--legacy-help		✓				
--version	✓	✓	✓	✓	✓	✓
--site-config-file	✓	✓	✓	✓	✓	✓
<i>filenames</i>	✓	✓	✓	✓	✓	✓
--xargs	✓	✓	✓	✓	✓	✓
--print-filenames	✓	✓	✓	✓	✓	✓
--copy-input		✓	✓	✓		✓
--pmap-file	✓	✓		✓	✓	✓
--plugin	✓	✓		✓	✓	✓
--python-file	✓	✓		✓	✓	✓
--output-path		✓	✓	✓	✓	✓
--no-titles		✓	✓	✓		✓
--no-columns		✓	✓	✓		✓
--column-separator		✓	✓	✓		✓
--no-final-delimiter		✓	✓	✓		✓
--delimited		✓	✓	✓		✓
--ipv6-policy		✓		✓		✓
--ip-format		✓		✓		✓
--timestamp-format		✓	✓	✓		✓
--integer-sensors		✓		✓		✓
--integer-tcp-flags		✓		✓		✓
--pager		✓	✓	✓		✓
--note-add	✓				✓	
--note-file-add	✓				✓	
--dry-run	✓			✓		



Table C.18: Parameters Common to Several Commands

Parameter	Description
<code>--help</code>	Prints usage description and exits
<code>--legacy-help</code>	Prints help for legacy switches
<code>--version</code>	Prints this program's version and installation parameters
<code>--site-config-file</code>	Specifies the name of the SiLK configuration file to use instead of the file in the root directory of the repository
<i>filenames</i>	Specifies one or multiple filenames as non-option arguments
<code>--xargs</code>	Specifies the name of a file (or <code>stdin</code> if omitted) from which to read input filenames
<code>--print-filenames</code>	Displays input filenames on <code>stderr</code> as each file is opened
<code>--copy-input</code>	Specifies the file or pipe to receive a copy of the input records
<code>--pmap-file</code>	Specifies a prefix-map filename and a map name as <i>mapname:path</i> to create a many-to-one mapping of field values to labels. For <code>rwfilter</code> , this creates new partitioning options: <code>--pmap-src-mapname</code> , <code>--pmap-dst-mapname</code> , and <code>--pmap-any-mapname</code> . For other tools, it creates new fields <code>src-mapname</code> and <code>dst-mapname</code> (see Section 6.2.7)
<code>--plugin</code>	For <code>rwfilter</code> , creates new switches and partitioning options with a plug-in program written in the C language. For other tools, creates new fields
<code>--python-file</code>	For <code>rwfilter</code> , creates new switches and partitioning options with a plug-in program written in Python. For other tools, creates new fields
<code>--output-path</code>	Specifies the output file's path
<code>--no-titles</code>	Doesn't print column headings
<code>--no-columns</code>	Doesn't align neat columns. Deletes leading spaces from each column
<code>--column-separator</code>	Specifies the character displayed after each column value
<code>--no-final-delimiter</code>	Doesn't display a column separator after the last column
<code>--delimited</code>	Combines <code>--no-columns</code> , <code>--no-final-delimiter</code> , and, if a character is specified, <code>--column-separator</code>
<code>--ipv6-policy</code>	Determines how IPv4 and IPv6 flows are handled when SiLK has been installed with IPv6 support (see Table C.17)
<code>--ip-format</code>	Chooses the format of IP addresses in output (see Table C.19)
<code>--timestamp-format</code>	Chooses the format and/or timezone of timestamps in output (see Table C.20)
<code>--integer-sensors</code>	Displays sensors as integers, not names
<code>--integer-tcp-flags</code>	Displays TCP flags as integers, not strings
<code>--pager</code>	Specifies the program used to display output one screen at a time
<code>--note-add</code>	Adds a note, specified in this option, to the output file's header
<code>--note-file-add</code>	Adds a note from the contents of the specified file to the output file's header
<code>--dry-run</code>	Checks parameters for legality without actually processing data

Table C.19: `--ip-format` Values

Value	Description
<code>canonical</code>	Displays IPv4 addresses as dotted decimal quad and most IPv6 addresses as colon-separated hexadectets. IPv4-compatible and IPv4-mapped IPv6 addresses will be displayed in a combination of hexadecimal and decimal. For both IPv4 and IPv6, leading zeroes will be suppressed in octets and hexadectets. Double-colon compaction of IPv6 addresses will be performed.
<code>zero-padded</code>	Octets are zero-padded to three digits, and hexadectets are zero-padded to four digits. Double-colon compaction is not performed, which simplifies sorting addresses as text.
<code>decimal</code>	Displays an IP address as a single, large decimal integer.
<code>hexadecimal</code>	Displays an IP address as a single, large hexadecimal integer.
<code>force-ipv6</code>	Display all addresses as IPv6 addresses, using only hexadecimal. IPv4 addresses are mapped to the <code>::FFFF:0:0/96</code> IPv4-mapped netblock.

Table C.20: `--timestamp-format` *format*, *modifier*, and *timezone* Values

Value	Description
<code>default</code>	Formats timestamps as <code>YYYY/MM/DDThh:mm:ss.sss</code> ( <code>rwcut</code> and <code>rwcount</code> may display milliseconds; <code>rwuniq</code> and <code>rwstats</code> never do)
<code>iso</code>	Formats timestamps as <code>YYYY-MM-DD hh:mm:ss.sss</code>
<code>m/d/y</code>	Formats timestamps as <code>MM/DD/YYYY hh:mm:ss.sss</code>
<code>epoch</code>	Formats timestamps as the number of seconds since 1970/01/01 00:00:00 UTC (UNIX epoch) <code>sssssssss.sss</code>
<code>no-msec</code>	Truncates milliseconds ( <code>.sss</code> ) from <code>sTime</code> , <code>eTime</code> , and <code>dur</code> fields – <code>rwcut</code> only
<code>utc</code>	Specifies timezone to use Coordinated Universal Time (UTC)
<code>local</code>	Specifies timezone to use the TZ environment variable or the system timezone

## Appendix D

# Additional Information on SiLK

*Network Traffic Analysis with SiLK* has been designed to provide an overview of data analysis with SiLK on an enterprise network. This overview has included the definition of network flow data, the collection of that data on the enterprise network, and the analysis of that data using the SiLK tool suite. The last chapter provided a discussion on how to extend the SiLK tool suite to support additional analyses.

This handbook provides a large group of analyses in the examples, but these examples are only a small part of the set of analyses that SiLK can support. The SiLK community continues to develop new analytical approaches and provide new insights into how analysis should be done. The authors wish the readers of this handbook good fortune in participation as part of this community.

### D.1 SiLK Support and Documentation

The SiLK tool suite is available in open-source form at <https://tools.netsa.cert.org/silk/>.

Before asking others to help with SiLK questions, it is wise to look first for answers in these resources:

***SiLK\_tool --help***: All SiLK tools (e.g., `rwfilter` or `rwcut`) support the `--help` option to display terse information about the syntax and usage of the tool.

**man pages**: All SiLK tools have online documentation known as manual pages, or `man` pages, that describe the tool more thoroughly than the `--help` text. The description is not a tutorial, however. `man` pages can be accessed with the `man` command on a system that has SiLK installed or via web links listed on the SiLK Documentation webpage at <https://tools.netsa.cert.org/silk/docs.html#manuals>.

***The SiLK Reference Guide***: This guide contains the entire collection of `man` pages for all the SiLK tools in one document. It is provided at <https://tools.netsa.cert.org/silk/reference-guide.html> in HTML format and at <https://tools.netsa.cert.org/silk/reference-guide.pdf> in Adobe® Portable Document Format (PDF) .

**SiLK Tool Suite Quick Reference booklet**: This very compact booklet (located at <https://tools.netsa.cert.org/silk/silk-quickref.pdf>) describes the dozen most used SiLK commands in a small (5.5" × 8.5") format. It also includes tables of flow record fields and transport layer protocols.

**SiLK FAQ:** This webpage answers frequently asked questions about the SiLK analysis tool suite. Find it at <https://tools.netsa.cert.org/silk/faq.html>.

**SiLK Tooltips:** This wiki contains tips and tricks posted by clever SiLK users. This useful information often is not immediately obvious from the standard documentation. Find it at <https://tools.netsa.cert.org/tooltips.html>.

**PySiLK Tooltips:** Tips for writing PySiLK scripts are provided at <https://tools.netsa.cert.org/confluence/display/tt/Writing+PySiLK+scripts>.

## D.2 FloCon Conference and Social Media

The CERT Division of the SEI supports FloCon<sup>®</sup>, an annual international conference devoted to large-scale data analysis for improving the security of networked systems—including flow analysis. More information on FloCon is provided at <http://www.flocon.org>.

Since the FloCon conference covers a range of network security topics, including network flow analysis, the conference organizers encourage ongoing discussions. In support of this, the following social networking opportunities are offered:

**@FloCon\_News Twitter account:** The FloCon conference organizers post notices related to FloCon here. View (and follow!) the FloCon tweets at [https://twitter.com/FloCon\\_News](https://twitter.com/FloCon_News).

**“FloCon Conference” LinkedIn member:** The FloCon Conference member page (<https://www.linkedin.com/in/flocon>) displays postings from the conference organizers, as well as from participants.

**“FloCon” LinkedIn group:** You can request membership to this private LinkedIn group at <https://www.linkedin.com/groups?gid=3636774>. Here, members discuss matters related to the FloCon conference and network flow analysis.

## D.3 Email Addresses and Mailing Lists

The primary SiLK email addresses and lists are described below:

**[netsa-tools-discuss@cert.org](mailto:netsa-tools-discuss@cert.org):** This distribution list is for discussion of tools produced by the CERT/CC for network situational awareness, as well as for discussion of flow usage and analytics in general. The discussion might be about interesting usage of the tools or proposals to enhance them. You can subscribe to this list at <https://lists.sei.cmu.edu>.

**[netsa-help@cert.org](mailto:netsa-help@cert.org):** This email address is for bug reports and general inquiries related to SiLK, especially support with deployment and features of the tools. It provides relatively quick response from CERT/CC users and maintainers of the SiLK tool suite. While a specific response time cannot be guaranteed, this address has proved to be a valuable asset for bugs and usage issues.

**[netsa-contact@cert.org](mailto:netsa-contact@cert.org):** This email address provides an avenue for recipients of CERT Situational Awareness Group technical reports to reach the reports’ authors. The focus is on analytical techniques. Public reports are provided at <https://www.cert.org/netsa/publications/>.

**[flocontact@cert.org](mailto:flocontact@cert.org)**: General email address for inquiries about FloCon.

**[flocommunity@cert.org](mailto:flocommunity@cert.org)**: This distribution list addresses a community of analysts built on the core of the FloCon conference. The list is not focused exclusively on FloCon itself, although it will include announcements of FloCon events.



# Appendix E

## Further Reading and Resources

This chapter gives helpful background information for many of the topics discussed in this guide. It is intended to be a starting point for further learning!

### E.1 Network Flow and Related Topics

Much has been written on NetFlow, SiLK, and related analysis. Here we provide a non-comprehensive list of examples you may wish to consider. Several related topics will enhance your ability to use SiLK effectively.

#### E.1.1 Technical Papers

Rick Hofstede, et al. *Flow Monitoring Explained: From Packet Capture to Data Analysis With NetFlow and IPFIX*, IEEE Communications Surveys and Tutorials (Volume 16, Issue 4, Fourth quarter 2014) <https://ieeexplore.ieee.org/document/6814316/?arnumber=6814316&tag=1>

T. Taylor, S. Brooks, J. McHugh. “NetBytes Viewer: An Entity-Based NetFlow Visualization Utility for Identifying Intrusive Behavior.” In: Goodall J.R., Conti G., Ma KL. (eds) *VizSEC 2007. Mathematics and Visualization*. Springer, Berlin, Heidelberg. [https://link.springer.com/chapter/10.1007/978-3-540-78243-8\\_7#citeas](https://link.springer.com/chapter/10.1007/978-3-540-78243-8_7#citeas)

T. Taylor, D. Paterson, J. Glanfield, C. Gates, S. Brooks and J. McHugh, “FloVis: Flow Visualization System,” *2009 Cybersecurity Applications and Technology Conference for Homeland Security*, Washington, DC, 2009, pp. 186-198. doi: 10.1109/CATCH.2009.18 <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=4804443&isnumber=4804414>

Jeff Janies, Red Jack. *Protographs: Graph-Based Approach to NetFlow Analysis*. FloCon 2011 [https://resources.sei.cmu.edu/asset\\_files/Presentation/2011\\_017\\_101\\_50576.pdf](https://resources.sei.cmu.edu/asset_files/Presentation/2011_017_101_50576.pdf)

M. Thomas, L. Metcalf, J. Spring, P. Krystosek and K. Prevost, “SiLK: A Tool Suite for Unsampled Network Flow Analysis at Scale,” *2014 IEEE International Congress on Big Data*, Anchorage, AK, 2014, pp. 184-191. doi: 10.1109/BigData.Congress.2014.34 <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=6906777&isnumber=6906742>

V. Marinov, J. Schoenwaelder. “Design of an IP Flow Record Query Language.” In: D. Hausheer, J.

Schoenwaelder (eds) *Resilient Networks and Services*. AIMS 2008. Lecture Notes in Computer Science, vol 5127. Springer, Berlin, Heidelberg, 2008

M. M. Najafabadi, T. M. Khoshgoftaar, C. Calvert and C. Kemp, “Detection of SSH Brute Force Attacks Using Aggregated Netflow Data,” *2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA)*, Miami, FL, 2015, pp. 283-288. doi: 10.1109/ICMLA.2015.20 <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7424322&isnumber=7424247>

Udaya Wijesinghe, Udaya Tupakula, Vijay Varadharajan. “An Enhanced Model for Network Flow Based Botnet Detection.” *Proceedings of the 38th Australasian Computer Science Conference (ACSC 2015)*, Sydney, Australia, 27 - 30 January 2015 <http://crpit.com/confpapers/CRPITV159Wijesinghe.pdf>

### E.1.2 Books on Network Flow and Network Security

Michael W. Lucas. *Network Flow Analysis*. no starch press. June 2010, ISBN-13: 978-1-59327-203-6 <https://nostarch.com/networkflow>

Omar Santos, *Network Security with NetFlow and IPFIX*. 2016 Cisco Systems, Inc, Cisco Press, Indianapolis, IN <http://www.ciscopress.com/store/network-security-with-netflow-and-ipfix-big-data-analytics-9781587144387>

## E.2 Bash Scripting Resources

Throughout this book, we use bash scripts to organize and execute collections of SiLK commands. There are many books, online classes, and web tutorials on Bash and its uses. Here are some that may be helpful.

### E.2.1 Online Tutorial

Online Shell Scripting Tutorial: <https://www.shellscript.sh>

### E.2.2 Books on Bash Scripting

Many books have been written on Bash and shell scripting in general. See the following publishers for their current list of titles relating to Bash.

**Google search for Bash scripting books:** <https://www.google.com/search?q=list+of+books+on+Bash+scripting>

#### O’Reilly

<https://www.oreilly.com>

<https://ssearch.oreilly.com/?q=bash>

#### No Starch Press, Inc.

<https://nostarch.com>

<https://nostarch.com/search/node/bash%20scripting>

#### John Wiley and Sons/WROX

[http://www.wrox.com/WileyCDA/Section/id-WROX\\_SEARCH\\_RESULT.html?query=bash](http://www.wrox.com/WileyCDA/Section/id-WROX_SEARCH_RESULT.html?query=bash)



**Addison Wesley**

<https://openlibrary.org/search?q=bash+scripting&mode=everything>

**Linux Training Academy**

<https://www.linuxtrainingacademy.com/books>

## E.3 Visualization

The following tools are useful for visually displaying SiLK data.

### E.3.1 Rayon

Rayon is a Python library and set of tools for generating basic, two-dimensional statistical visualizations. Rayon can be used to automate reporting; provide data visualization in command-line, GUI or web applications; or do ad-hoc exploratory data analysis.

<https://tools.netsa.cert.org/rayon/index.html>

### E.3.2 FloViz

FloViz is a comprehensive and extensible set of visualization tools. It is integrated with the SiLK tool suite via a relational database that stores data such as sets and multisets (bags) that are derived from NetFlow and similar sources.

<https://web.cs.dal.ca/~sbrooks/projects/NetworkVis/index.html>

### E.3.3 Graphviz - Graph Visualization Software

Network flow records can be visualized as directed graphs. Graphviz can be used to produce such visualizations.

<https://www.graphviz.org>

### E.3.4 The Spinning Cube of Potential Doom

First implemented as a demonstraton, it is an interesting display of network traffic at the Supercomputing conference in 2004. It has since taken on a life of its own.

<http://www.nersc.gov/news-publications/nersc-news/nersc-center-news/1998/cube-of-doom>

Stephen Lau. “The Spinning Cube of Potential Doom.” *Commun. ACM* 47, 6 (June 2004), 25-26. <https://dx.doi.org/10.1145/990680.990699>

## E.4 Networking Standards

*Service Name and Transport Protocol Port Number Registry*  
<https://www.iana.org/assignments/service-names-port-numbers>

*RFC 871, A Perspective on the ARPANET Reference Model.*  
<https://tools.ietf.org/html/rfc871>

*ETF RFC 3954, Cisco Systems NetFlow Services Export Version 9, 2004*  
<https://www.ietf.org/rfc/rfc3954.txt>

*RFC 4291, IP Version 6 Addressing Architecture*  
<https://tools.ietf.org/html/rfc4291>

*RFC 6890, Special-Purpose IP Address Registries*  
<https://tools.ietf.org/html/rfc6890>.

IPFIX standards:

RFC3917: Requirements for IP Flow Information Export (IPFIX)

RFC3955: Candidate Protocols for IP Flow Information Export (IPFIX)

RFC5101: Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of IP Traffic Flow Information (IPFIX)

RFC5102: Information Model for IP Flow Information Export

RFC5103: Bidirectional Flow Export Using IP Flow Information Export

RFC5153: IPFIX Implementation Guidelines

RFC5470: Architecture for IP Flow Information Export

RFC5471: Guidelines for IP Flow Information Export (IPFIX) Testing

RFC5472: IP Flow Information Export (IPFIX) Applicability

RFC5473: Reducing Redundancy in IP Flow Information Export (IPFIX) and Packet Sampling (PSAMP) Reports

# Index

- ACK, [58](#), [169](#)
- address relationships, [53](#)
- address types, [188](#)
- addressing, [164](#)
- advanced analysis, *see* [exploratory analysis](#)
- all, [7](#)
- anonymization, [120](#)
- application layer, [162](#)
  
- bags, [67](#), [137](#)
  - adding, [110](#)
  - binning, [72](#)
  - comparing contents of, [73](#)
  - counters, [67](#), [68](#)
  - cover sets, [112](#)
  - creating, [68](#), [69](#)
  - displaying counts and key values, [71](#)
  - dividing, [110](#), [112](#)
  - finding scanners, [114](#)
  - formatting display of, [73](#)
  - intersecting with IPsets, [74](#)
  - key values, [67](#), [68](#)
  - multiple, [68](#)
  - multiplying by a scalar value, [110](#), [112](#)
  - relationship to IPsets, [67](#), [112](#)
  - relationship to prefix maps, [121](#)
  - sensor inventories, [93](#)
  - subtracting, [110](#)
  - summarizing NTP traffic, [104](#)
  - summarizing web traffic, [69](#)
  - thresholding, [71](#), [112](#), [114](#)
  - viewing file information, [90](#)
- Bash, [142](#), [179](#)
- basic analysis, *see* [single-path analysis](#)
- behavioral analysis, [35](#)
- bins, [28](#), [30](#)
  - allocating flows, packets and bytes, [63](#)
  - bag counts, [72](#)
  - bottom-*N* and top-*N*, [30](#)
  - counting traffic volumes, [32](#)
  - plotting, [33](#)
  - size of, [29](#)
  - skipping zero-size, [35](#)
  - varying sizes of, [34](#)
- bottom-*N* lists, [30](#)
- bytes, [3](#), [17](#)
  - counting, [28](#), [32](#)
  - destination port usage, [47](#)
  - field number, [3](#)
  - filtering by byte count, [29](#), [44](#)
  - sorting by count, [35](#)
  - thresholding, [64](#)
  
- case studies, [43](#), [93](#), [131](#)
  - dataset for, [13](#)
  - scripting, [96](#)
- CIDR notation, [37](#), [47](#), [122](#), [164](#)
- client, [58](#), [61](#)
- command line, [173](#)
- complement intersect, [74](#)
- conditional fields, [152](#)
- counters, [67](#), [71](#)
  - comparing, [73](#)
  - examples of, [68](#)
- country codes, [3](#), [125](#), [188](#)
- cover sets, [74](#), [112](#), [114](#)
  
- data structures, [148](#)
- dataset for examples, [13](#)
- date format, [8](#)
- destination IP address, [3](#)
  - CIDR notation for, [47](#)
  - creating IPsets, [38](#)
  - displaying, [25](#)
  - sorting by, [35](#)
- destination IP type, [3](#)
- destination port, [3](#), [47](#)
  - displaying, [25](#)
- DHCP, [31](#)
- dIP, *see also* [destination IP address](#), [3](#)
  - displaying, [25](#)
  - field number, [3](#)

- distributed environments, 147
- DNS, 31, 40, 47, 66, 77, 78, 171
- Domain Name System, *see* DNS
- domain names, 40
- dPort, *see also* destination port, 3
  - displaying, 25
  - field number, 3
- dType, 3
  - field number, 3
- duration, 3
  - field number, 3
  - filtering by, 44
- Dynamic Host Configuration Protocol, *see* DHCP31
- end time, 3
- enterprise network, 6, 93, 162, 165
  - sensor information, 6, 17
  - subnet, 39
- entry, 73
- eTime, *see also* end time, 3
  - field number, 3
- exploratory analysis, 10, 99, 106
  - case study, 131
  - commands, 107
  - dataset for examples, 13
  - example of, 101
  - overview of, 99
  - relationship to single and multi-path analysis, 99
  - starting points, 101
  - workflow, 100
- ext2ext, 7
- FCCX dataset, 13, 122
  - network diagram for, 131
- fields, 3
  - character string, 154
  - conditional, 152
  - extending with PySiLK, 151
  - names and numbers for, 3, 26, 192
  - sorting by, 35
  - stored vs. derived, 3
- FIFO, 53, 178
- files
  - appending, 114
  - bag, 90
  - binary, 24
  - combining, 114
  - filtering, 23
  - flow repository, 7, 18
  - IPset, 37, 90
  - network flow record, 3, 23, 26, 29
  - prefix map, 90, 122
  - simple anonymization, 120
  - splitting, 116
  - temporary, 36
  - variable record length, 27
  - viewing contents of, 24
  - viewing information about, 26, 90
- filtering, 11, 44, 53
  - all destinations, 56
  - by byte count, 29
  - by country code, 125
  - by prefix value, 125
  - complex, 56, 148
  - extending with PySiLK, 143
  - inbound client traffic, 61
  - inbound server traffic, 58, 59
  - inbound TCP traffic, 58
  - internal, external, and non-routable addresses, 125
  - IPsets, 66
  - isolating behaviors of interest, 66
  - low-packet flows, 61
  - manifold, 57
  - outbound client traffic, 61
  - outbound server traffic, 61
  - overlapping traffic, 57
  - pass-fail, 56, 57
  - prefix maps, 125
  - removing unwanted flows, 47, 59
  - role of partitioning parameters, 21
  - tuple files, 107
- FIN, 169
- five-tuple, 3, 107
- flags, *see* TCP flags
- FloCon conference, 220
- FloViz, 225
- flow data, *see* network flow records
- flow file, *see* network flow records
- flow label, 2
- flow record, *see* network flow records
- flow repository, 2, 7
  - querying, 16
  - retrieving records from, 21
  - structure of, 7
  - viewing time information, 20
- flow type, 3
  - field number, 3
  - retrieving and filtering data, 23

- viewing, 18
- flows, 1, 4
  - approximating over time, 63
  - counting, 28, 30, 32
  - destination port usage, 47
  - grouping, 82
  - low-byte vs. high-byte, 28
  - mismatched, 135
  - packets, 169
  - split, 4
  - statistical summaries of, 28
  - thresholding, 64
- formulate, 10
- Graphviz, 225
- groups, 82
  - by prefix value, 126
  - creating from IPsets, 86
  - matching queries and responses, 86
  - sorting, 83
  - summarizing, 84
  - thresholding, 84
- help with SiLK, 181
- here-documents, 177
- ICMP, 31, 49, 171
- in, 7, 58, 70
- inicmp, 7
- innull, 7
- input parameters, 21
- int2int, 7
- intermediate analysis, *see* multi-path analysis
- intrusion detection signatures, 130
- inweb, 7, 70
- IP, 161, 162, 164
- IP addresses, 2, 70, 164, 187
  - associating with sensor, 94
  - bags, 67
  - binning, 72
  - CIDR notation, 37, 122
  - counting, 112
  - displaying prefix values, 128
  - filtering by CIDR block, 47
  - format of, 164
  - IPsets, 37
  - IPv4 vs. IPv6, 25, 164
  - labeling with prefix maps, 122
  - limiting, 74, 77
  - masking, 75
  - removing from bags, 112
  - reserved, 165
  - resolving to domain name, 40
  - thresholding, 64
  - wildcard notation, 38
- IPFIX, 3, 226
- IPsets, 17, 37, 40, 76
  - algebraic operations, 77
  - by sensor, 78
  - combining, 77, 96
  - counting members of, 38, 80
  - cover set, 112
  - creating, 37, 76, 112
  - creating bags from, 69
  - difference between, 77, 112
  - displaying members of, 38, 80
  - extracting from bags, 74, 112
  - filtering with, 66
  - finding scanners, 112
  - generating from `rwfilter`, 66
  - generating from bags, 114
  - grouping, 86
  - intersecting, 78
  - intersecting with bags, 74
  - limiting IP addresses, 77
  - members of, 78
  - relationship to bags, 67, 112
  - relationship to prefix maps, 121
  - sensor inventories, 93
  - summary statistics for, 39
  - symmetric difference, 78
  - time period for, 77
  - viewing file information, 90
- IPv4, 25, 110, 163, 187, 213
  - address format, 164
  - prefix maps, 122
  - reserved addresses, 165
- IPv6, 25, 163
  - address format, 164
  - prefix maps, 122
  - reserved addresses, 165
- iterating, 12, 55, 99
- key values, 67, 71
  - comparing, 73
  - examples of, 68
  - in cover sets, 74
- load scheme, 63, 191

- manifold, 53, 57–59
  - command examples, 58, 59, 62
  - filtering low-packet flows, 61
  - non-overlapping, 57
  - overlapping, 57
  - use in profiling, 66
- matched groups, 86
- matching
  - flows, 135
  - incomplete sessions, 135
  - sorting before, 87
- message, 162
- mkfifo, 53
- multi-path analysis, 9, 51
  - case studies, 93
  - common commands, 56
  - dataset for examples, 13
  - exploring relationships, 53
  - interpreting results, 55
  - overview of, 51
  - pitfalls, 55
  - relationship to exploratory analysis, 99
  - relationship to single-path analysis, 53
  - scripting, 96
  - workflow, 51
- named pipes, 53, 178
- NetFlow, 2
- network application, 162
- network flow, *see* flows, network flow records
- network flow records, 1–3
  - collection of, 6
  - combining files, 114
  - counting by prefix value, 127
  - counting in file, 26
  - creating bags from, 68
  - creating from text, 119
  - creating IPsets from, 37
  - fields in, 3
  - filtering, 21
  - flow label, 2
  - generation, 4
  - grouping, 82
  - labeling with prefix maps, 121
  - pulling from repository, 21
  - querying, 16
  - removing unwanted flows, 47, 59
  - sorting, 35
  - splitting files, 116
  - storage of, 7
  - thresholding, 64
  - time and date, 8
  - variable length, 27
  - viewing in text format, 24
  - where collected, 6
- network layer, 162
- network mask, 40, 75
- Network Time Protocol, *see* NTP
- network traffic, 2, 7
  - anomalies, 133
  - asymmetric and missing data, 95
  - categorizing, 57
  - counting, 28, 32
  - excluding, 17
  - filtering, 21
  - finding commonly-used protocols, 30
  - plotting, 33
  - profiling around an event, 43
  - profiling with IPsets, 78
  - summarizing, 28
  - types of, 7
  - web services, 52
- network traffic analysis, 9
  - case studies, 43, 93, 131
  - exploratory, 10, 99
  - multi-path, 9, 51
  - single-path, 9, 15
  - workflow, 10
- next-hop IP, 3, 135
  - use in groups, 84
  - use in matching, 89
  - use in sets, 86
- nhIP, *see also* next-hop IP, 3
  - field number, 3
- NTP, 101
- Open Systems Interconnection model
  - see* OSI, 161
- OSI, 161
- other, 7
- out, 7
- outicmp, 7
- outnull, 7
- output parameters, 183
  - for `rwfilter`, 22
  - list of, 188
- outweb, 7
- packets, 3, 17, 162, 169
  - and TCP/IP, 161

- counting, [28](#), [32](#)
- field number, [3](#)
- TCP, [168](#)
- thresholding, [64](#)
- time distribution, [63](#)
- UDP and ICMP, [172](#)
- PAGER environment variable, [26](#)
- partitioning, [10](#)
- partitioning parameters, [21](#), [183](#)
  - data structures for, [148](#)
  - extending with PySiLK, [143](#), [145](#)
  - list of, [184](#)
  - Python boolean expressions, [147](#)
- pass-fail filtering, [21](#), [57](#)
- physical layer, [163](#)
- pipes, [21](#), [22](#), [57](#), [114](#), [133](#), [177](#)
  - named, [53](#)
  - use with `rwcut`, [24](#)
  - use with `rwsort`, [35](#)
  - use with `rwuniq`, [29](#)
  - when to use, [29](#)
- plots, [33](#)
- plug-ins, [141](#), [142](#)
  - code examples, [146](#), [152](#), [154](#), [156](#), [158](#)
  - use with `rwfilter`, [143](#)
  - with `silkpython`, [142](#)
- port knocking, [145](#)
- port-protocol pairs, *see* protocol-port pairs
- ports, [2](#)
  - summarizing traffic with bags, [67](#)
  - traffic anomalies, [133](#)
- prefix maps, [121](#)
  - counting by prefix value, [126](#)
  - country codes, [125](#)
  - creating, [122](#)
  - displaying prefix values, [126](#)
  - filtering with, [125](#)
  - grouping by prefix value, [126](#)
  - internal, external, and non-routable addresses, [125](#)
  - naming, [122](#)
  - querying, [128](#)
  - relationship to bags and sets, [121](#)
  - sorting by prefix value, [126](#)
  - statistics by prefix value, [126](#)
  - user-defined vs. predefined, [121](#), [125](#)
  - viewing file information, [90](#)
- presentation layer, [162](#)
- protocol, [3](#), [161](#), [162](#), [168](#)
  - behavioral analysis of activity, [35](#)
  - displaying, [25](#)
  - field number, [3](#)
  - ICMP, [31](#)
  - sorting by, [35](#)
  - summarizing traffic with bags, [67](#)
  - TCP, [31](#)
  - UDP, [31](#)
- protocol layers, [168](#)
- protocol-port pairs, [121](#), [122](#)
  - building prefix map, [130](#)
  - displaying prefix values, [128](#)
  - labeling with prefix maps, [122](#)
  - notation for, [122](#)
- PSH, [169](#)
- PySiLK, [141](#), [142](#)
  - code examples, [145–149](#), [152](#), [154](#), [156](#), [158](#)
  - complex filtering with, [148](#)
  - conditional values, [152](#)
  - defining character string fields, [154](#)
  - defining key and summary value fields, [158](#)
  - distributed environments, [147](#)
  - extending fields, [151](#)
  - preserving state information, [143](#), [145](#)
  - programming with, [142](#)
  - requirements, [142](#)
  - use with `rwcut`, [152](#)
  - use with `rwfilter`, [143](#)
  - use with `rwsort`, [152](#)
- Python, [119](#), [141](#), [142](#)
  - boolean expressions and `rwfilter`, [147](#)
  - data structures as partitioning parameters, [148](#)
- PYTHONPATH environment variable, [142](#)
- queries, [10](#), [16](#), [21](#), [44](#)
  - complex, [56](#)
  - in exploratory analysis, [101](#)
  - narrowing focus of, [24](#)
  - pass-fail, [21](#)
  - prefix maps, [128](#)
- query matching, [87](#)
- Rayon, [225](#)
- repository, *see* flow repository
- reserved IP addresses, [165](#)
- response matching, [87](#)
- reverse DNS lookup, [40](#)
- RIP, [31](#)
- routers, [3](#)
- routing, [164](#)

- Routing Information Protocol, *see* RIP
- RST, 169
- rwappend, 114, 198
  - command examples, 107
- rwbag, 68, 205
  - command examples, 53, 68, 69, 94, 97, 114, 137
- rwbagbuild, 69, 206
  - command examples, 70, 104
  - file format, 70
- rwbagcat, 208
  - command examples, 53, 70–75, 104, 110, 112
  - dividing and multiplying bags, 112
- rwbagtool, 73, 209
  - adding and subtracting bags, 110
  - command examples, 74, 75, 94, 97, 110, 112, 114, 137
  - dividing and multiplying bags, 110
  - extracting cover sets, 74
  - intersecting bags and IPsets, 74
  - logical operations on key/counter values, 73
- rwcat, 114, 197
  - command examples, 116
- rwcount, 17, 32, 191
  - command examples, 33, 103
  - default bin size, 34
  - load scheme, 63, 191
  - prefix maps, 126
  - skip zero size bins, 35
  - time series plots, 33
- rwcut, 17, 24, 192
  - command examples, 25, 35, 47, 75, 84, 86, 87, 89, 120, 126, 127, 152, 154, 156
  - conditional values, 152
  - default fields, 26
  - defining character string fields, 154
  - delimiters, 26
  - display order of fields, 26
  - displaying fields, 25
  - extending with PySiLK, 152
  - number of records, 25
  - prefix maps, 126
  - use in behavioral analysis, 35
- rwfileinfo, 26, 90, 210
  - command examples, 27, 90, 116, 120, 126
  - default fields, 27
- rwfilter, 8, 17, 21, 28, 183
  - code examples, 147, 148
  - command examples, 22, 23, 30, 37, 45, 47, 53, 58, 59, 62, 64–70, 75, 77, 78, 86, 87, 94, 97, 102–104, 107, 108, 112, 114, 116, 119, 120, 126, 133, 136–138
  - complex filtering, 56
  - data flow in, 22
  - displaying file names, 24
  - extending with PySiLK, 143, 148
  - filtering by byte count, 29
  - finding low-packet flows, 61
  - IDS signatures, 130
  - in distributed environments, 147
  - input parameters, 21
  - IP address, 23
  - manifold, 57–59, 61
  - miscellaneous parameters, 189
  - multiple input files, 114
  - output parameters, 22, 188
  - parameter relationships, 22
  - partitioning parameters, 21, 184
  - pass-fail filtering, 22, 23, 56
  - plug-ins, 146
  - prefix maps, 125
  - Python boolean expressions, 147
  - selection parameters, 21, 183
    - use as partitioning parameters, 23
  - sensor, 23
  - start and end times, 22
  - tuple files, 107
  - type, 23
  - use in behavioral analysis, 35
  - use with multiple files, 24
- rwgroup, 83, 214
  - command examples, 84, 86
  - defining character string fields, 154
  - extending with PySiLK, 154
  - grouping by session, 83
  - prefix maps, 126
  - sorting before use, 83
  - thresholding, 84
- rwidsquery, 130
- rwmatch, 86, 213
  - command examples, 87, 89, 136
  - sorting before use, 83, 87
- rwnetmask, 75, 196
  - command examples, 75
- rwpmabuild, 122, 211
  - command examples, 123, 130
  - input file format, 122
- rwpmlookup, 128, 212
  - command examples, 130



- `rwresolve`, 40
  - command example, 40
- `rwscan`, 70
  - command examples, 70
- `rwset`, 17, 37, 201
  - command examples, 37, 67, 77, 78, 86, 94, 97, 112
  - use with `rwfilter`, 66
- `rwsetbuild`, 17, 37, 76, 204
  - command examples, 38, 75, 76, 94, 97
- `rwsetcat`, 17, 38, 80, 202
  - command examples, 39, 40, 74, 77, 78, 80, 81, 94, 97, 112, 114
  - displaying network structure, 81
- `rwsetmember`, 78
  - command examples, 78
- `rwsettool`, 77, 203
  - command examples, 77, 78, 80, 94, 97, 112, 114
  - difference, 77
  - displaying repository dates, 77
  - intersecting IPsets, 78
  - symmetric difference, 78
- `rwsiteinfo`, 6, 17, 182
  - command examples, 18, 77, 78
  - displaying sensors, 18
  - displaying traffic information, 18
- `rwsort`, 35, 194
  - command examples, 35, 84, 87, 89, 127, 136, 154
  - conditional values, 152
  - defining character string fields, 154
  - extending with PySiLK, 154
  - field numbers, 36
  - multiple files, 36
  - multiple input files, 114
  - prefix maps, 126
  - sort order, 35, 36
  - sorting before `rwgroup` and `rwmatch`, 83, 87
  - use in behavioral analysis, 35
- `rwsplit`, 116, 199
  - command examples, 116, 119
- `rwstats`, 17, 28, 190
  - command examples, 31, 47, 66, 119, 136
  - compared to `rwuniq`, 32, 47
  - defining character string fields, 154
  - defining key and summary value fields, 158
  - extending with PySiLK, 154
  - prefix maps, 126
  - required fields, 31
  - thresholding on compound keys, 65
    - top-*N* and bottom-*N* lists, 30
    - top-*N* lists, 46
- `rwtuc`, 119, 200
  - command examples, 120
- `rwuniq`, 17, 28, 29, 195
  - command examples, 30, 45, 47, 64, 65, 67, 69, 102, 104, 108, 127, 133, 137, 138, 156, 158
  - compared to `rwstats`, 32, 47
  - compound keys, 65
  - defining character string fields, 154
  - defining key and summary value fields, 158
  - extending with PySiLK, 154
  - prefix maps, 126
  - profiling, 66
  - profiling traffic, 102
  - required parameters, 29
  - specifying ranges, 64
  - summarizing web traffic, 69
  - thresholding, 64
- scanning, 112, 114, 140
  - detecting with `rwscan`, 70
  - finding port scanners, 133
  - use of bags, 69
- scripting languages, 141, 179
  - use with `rwtuc`, 119
- selection parameters, 21, 183
  - list of, 183
  - use as partitioning parameters, 23
- sensor, 3
  - class, 8
  - displaying information for, 17
  - field number, 3
  - inventories, 93
  - locations, 6
  - network flow collection, 3
  - retrieving data for, 23
  - type, 8
- sensors, 78
- server, 58, 59, 61
- services, 171
- session layer, 162
- sessions
  - grouping, 83
  - matched groups, 86
- sets, *see* IPsets
- shell scripts, 96, 99, 142, 179
  - command examples, 97, 116, 138
  - examples of, 53
  - help with, 224

- SiLK, 1
  - analysis types, 9, 15, 51
  - applying workflow, 44
  - case studies, 43, 93, 131
  - common parameters, 215
  - email addresses, 220
  - flow repository, 7, 8
  - help with, 181, 219
  - IPv4 and IPv6 fields, 163
  - repository, *see* flow repository
  - source files, 219
  - tool suite, 8
  - use with Python, 142
  - version, 181
  - visualization tools, 225
  - wildcard notation for IP addresses, 38
  - workflow, 10
- SILK\_COUNTRY\_CODES environment variable, 125
- SILK\_IPV6\_POLICY environment variable, 25
- SILK\_PAGER environment variable, 26
- silkython, 141, 142
- Simple Network Management Protocol, *see* SNMP
- single-path analysis, 9, 15
  - behavioral analysis, 35
  - case studies, 43
  - common SiLK commands for, 17
  - dataset for examples, 13
  - interpreting, 49
  - overview of, 15
  - relationship to exploratory analysis, 99
  - relationship to multi-path analysis, 51
  - workflow, 15, 17
- sIP, *see also* source IP address, 3
  - displaying, 25
  - field number, 3
- SMTP, 171
- SNMP, 31
- sorting
  - by field number, 36
  - by prefix value, 126
  - multiple files, 36
  - on multi-field values, 152
  - order, 36
  - with `rwsort`, 35
  - with `rwuniq`, 29
- source IP address, 3
  - CIDR notation for, 47
  - creating IPsets, 38
  - displaying, 25
- source IP type, 3
- source port, 3
  - displaying, 25
- sPort, *see also* source port, 3
  - displaying, 25
  - field number, 3
- standard input, 38
- standard output, 29
- standards, 226
- start time, 3
  - displaying, 25
- state information, 143
- stderr, 176
- stdin, 176
- stdout, 176
- sTime, *see also* start time, 3, 29
  - field number, 3
- sType, 3
  - field number, 3
- subnet, 39, 74
  - in IPset, 39
- summary record, 84
- SYN, 58, 169
- TCP, 31, 161, 162, 168, 171
  - filtering with TCP flags, 58, 59, 61
  - finding suspicious requests, 133
  - header, 168
  - initial flags, 3
    - field number, 3
  - low-packet flows, 61
  - matching sessions, 89
  - packets, 169
  - selecting TCP flows, 53
  - session flags, 3
    - field number, 3
  - state machine, 169
  - summarizing web traffic, 69
- TCP flags, 17, 169, 187
  - client vs. server, 62
  - client-server communication, 58
  - filtering with, 58, 59, 61
- TCP/IP, 161, 168
  - protocol layers, 161
- text
  - converting to network flow records, 119
  - creating bags from, 69
  - creating IPsets from, 37, 76
  - creating prefix maps from, 122
  - viewing flow records as, 24

thresholding, [31](#), [32](#), [64](#), [84](#), [112](#), [114](#)  
    bag counts, [71](#)  
    compound keys, [65](#)  
    min-max, [64](#)  
    multiple parameters, [64](#)

time bins, *see* bins

time distribution of packets, [63](#)

time format, [8](#), [22](#)

time range, [188](#)

timing relationships, [55](#)

top-*N* lists, [30](#), [46](#)

traffic, *see* network traffic

Transmission Control Protocol, *see* TCP

transport layer, [162](#)

tuple files, [107](#)

tuple files:example of, [108](#)

type, *see* flow type, [7](#), [8](#), [29](#)

UDP, [31](#), [49](#), [66](#), [162](#), [171](#)  
    matching sessions, [89](#)  
    NTP, [101](#)

UNIX, [173](#)

URG, [169](#)

volume relationships, [55](#)

web services, [52](#), [171](#)

web traffic, [69](#)

wildcard notation for IP addresses, [38](#), [187](#)

workflow, [44](#)

yaf, [3](#)



Network Traffic Analysis with SiLK  
Analyst's Handbook for SiLK  
Version 3.12.0 and Later  
August 2018