



PCAP Anonymizer Quick Start Guide

Contents

Introduction	2
Basic Usage	2
Best Practices.....	3
Generate and Customize Options File	5
Anonymizer Options	5
Secret Keys.....	6
Anonymize MAC Addresses in Ethernet Frame Headers	6
Remove or Replace VLAN ID Values in 802.1Q Header	7
Recalculate Transport and IP Layer Checksums	7
Anonymize IP Addresses.....	8
Remove or Anonymize ARP Messages	8
Remove or Anonymize ICMP Messages	9
Packet Truncation to User-Specified Length	10
BPF Filtering of Packets Read	10
Payload Obfuscation	10
Network Graph Views of IP Communications	11

Introduction

Anonymizer is a command-line utility for anonymizing packet capture files. In this context, anonymization means protecting identifying information about network endpoints and their users.

Anonymizer strives to do a good job at protecting sensitive data while preserving the analytical value of the original PCAP. There is no one-size-fits-all approach, no perfect tool or technique, and different packet artifacts are important to different folks. So, another main design goal is to provide flexibility and control over how anonymization tasks are performed.

Anonymizer assumes sensible defaults and the only required argument is the source PCAP file. If you choose to omit all other flags, Anonymizer will perform full anonymization of MAC Addresses and IP Addresses found in the Ethernet and IP headers respectively. The resulting output file (the *generated* trace) is generally suitable for common PCAP sharing use cases. However, Anonymizer includes many other command-line flags dealing with PCAP data transformation, reporting and analysis.

Basic Usage

Using anonymizer in its default state is simple. Download it and run the executable, passing in the `-source` flag to specify the source PCAP file you wish to anonymize.

By default, without any configuration flags set, **Anonymizer will anonymize the IP addresses and MAC addresses found in every Ethernet Frame and IP header read from the source PCAP.** This is done by encrypting IP addresses and replacing MAC addresses with randomly generated values. IP and Transport header checksums will also be recalculated but no other changes are made to any frame or packet.

This default configuration covers many use cases, but Anonymizer offers many other features to make the outcome more concise and effective. To see the full list of supported features and additional usage information use the `-h` or `-help` flag. We also cover that in more detail below.

If you're looking to get started quickly, follow these simple steps:

1. Download the software binaries, documentation, and sample PCAPs:

Dynamite Analytics has set up a temporary webpage for the DOE PCAP Anonymizer project to allow the project reviewers to anonymously download the software and documentation. This page will be eventually replaced with the official software download that will require user registration.

PCAP Anonymizer project page: <https://dynamite.ai/anonymizer-project>

- Software binaries are available for Linux and macOS
- Sample PCAP files can be download from PacketTotal

2. Set executable permissions on the binary:

```
chmod +x anonymizer-linux-v0.5.0.bin
```

3. On macOS disable the verified app check:

```
sudo spctl --master-disable
```

4. Run the binary against a PCAP file:

```
./anonymizer -source my_sensitive.pcap
```

When complete, Anonymizer will generate a JSON formatted report on the console like the one below; and a new, anonymized, PCAP file in the directory specified with the `-outdir` flag or in `/tmp` by default.

```
{
  "source_file": "vlan.cap",
  "output_file": "/tmp/vlan.cap-anonymized",
  "packets_read": 395,
  "packets_analyzed": 214,
  "packets_written": 214,
  "packets_discarded": 0,
  "packets_ip_anonymized": 200,
  "frames_mac_anonymized": 381,
  "packets_vlan_anonymized": 0,
  "packets_vlan_removed": 0,
  "packets_truncated": 0,
  "fragments_truncated": 0,
  "payloads_overwritten": 0,
  "packet_write_errors": 0,
  "packet_modify_errors": 181,
  "packet_analysis_errors": 0,
  "packet_read_errors": 0,
  "serialize_layer_errors": 0,
  "unserializeable_layers": 0,
  "packet_serialize_errors": 0,
  "ip_layer_not_found": 161,
  "ethernet_layer_not_found": 0,
  "mapped_ip_addresses": 19,
  "mapped_mac_addresses": 57,
  "mapped_vlan_ids": 0,
  "arp_ignored": 4,
  "arp_removed": 0,
  "arp_anonymized": 0,
  "icmp_ignored": 10,
  "icmp_removed": 0,
  "icmp_anonymized": 0,
  "ip_encryption_errors": 0,
  "functions_run": [
    "packet.arpIgnore",
    "packet.icmpIgnore",
    "packet.macFullAnonymize",
    "packet.ipFullAnonymize"
  ]
}
```

Best Practices

We have begun to refine sound practices for using Anonymizer and sharing network trace files publicly in general. Preparation is key, and to prepare a PCAP file for sharing helps reduce risk of information

leakage but it does not eliminate all risk. There is also an inverse relationship between data anonymization and the analytical value – or the legitimacy and usefulness - of the generated trace. As more data obfuscation protections are applied – in most cases - the analytical value of the generated trace is reduced.

For these reasons, we recommend you be as concise as possible when selecting traffic to be shared, and only applying protections you deem absolutely necessary.

- **Ensure you have authority and permission to capture and share traffic** – You must ensure you have permission to capture and share network traces before doing so. Network traffic traces contain a trove of sensitive information, in various forms, and are owned by the organization within which they are captured. Only sharing what you have the authority and permission to share, helps limit your liability and reduce the risk of unintentional information leakage.
- **Do not reuse secret keys** – While convenient, the reuse of secret keys increases the potential for an adversary to successfully decrypt encrypted values. Using randomly generated keys helps ensure a given key is only useful in decrypting a single trace file.
- **Use a source BPF filter** – Use capture filters or Anonymizer’s BPF option to restrict packets read from the source trace file. This is one of the most effective ways to reduce the risk of accidental information leakage. By only including the traffic that is intended to be shared, no other traffic poses a risk.
- **Target IP Address Anonymization** – Applying IP anonymization to all IP addresses can have adverse effects on the fidelity of the generated trace file. Under most circumstances, it is sufficient to protect IP addresses under one’s ownership and control and leave others unmodified. This helps preserve more analytical value in the generated trace file while still protecting the identity of assets of interest.
- **Remove VLAN Headers** – Unless specifically needed, VLAN headers and ID’s they contain provide little value or context to researchers working with an anonymized trace file. While they in themselves do not pose significant identity exposure risk, adversaries can use VLAN information to map a victim’s network environment. Unless otherwise needed, we recommend you remove VLAN headers from anonymized packets to be shared using the `-vlan remove` flag or equivalent option.
- **Remove ICMP, ARP and other control messages** – Address resolution and control protocol messages pose elevated risk for privacy and asset identity leakage and should be discarded if not specifically required to satisfy a particular use case. By using restrictive BPF filters as described above Anonymizer can prune unwanted messages from a generated trace file, greatly reducing the risk of identity and sensitive information leakage.

This is not an exhaustive list, but a few recommendations that are applicable to most use-cases. Of course, anonymization needs can vary widely depending on the intended use of the anonymized trace and Anonymizer provides an extensive interface for customizing its anonymization features.

Generate and Customize Options File

Anonymizer provides many command-line flags as well as a customizable configuration file alternative, referred to as the *Options* file.

The easiest way to produce an Options file is to have Anonymizer do it for you by passing the `-print-options` flag. As with before, with no other flags set, Anonymizer will generate and print a YAML formatted options file like the one shown below - note the `-source` flag is *always* required:

```
./anonymizer -source ~/pcaps/vlan.pcap -print-options
data_link:
  mac_addresses: full
  vlan_ids: ignore
  arp_messages: ignore
network:
  ip_addresses: full
  ip_length: false
  icmp_messages: ignore
application:
  truncate_payload: false
  truncate_length: 1024
  mask_payload: false
key: b96ce879cdeeff578787bd796b00d7a5
store_mappings: false
gen_opts: false
source_file: /Users/jarvis/pcaps/vlan.pcap
output_directory: /tmp
output_file: vlan.pcap-anonymized
output_format: pcap-ng
cidrs: ""
analysis:
  graph: false
fix_checksums: true
pcapng_annotate: true
```

This output can simply be redirected to a file and edited in any text editor. Alternatively, you can specify the `-store-options` flag and Anonymizer will do the same thing, only store the Options file in YAML format in the specified output directory (or /tmp by default). For more information on each available option see the section below. When you are finished customizing the Options file, save it, and pass it as an argument to the `-options-file` flag as shown below:

```
./anonymizer -source ~/pcaps/vlan.pcap -options-file my_anon_options.yml
```

Anonymizer Options

The following section describes the Anonymizer features that can be enabled and configured via the Options file described above or by supplying the equivalent command-line flags.

Secret Keys

Anonymizer includes a simple secret key framework that serves as a basis for encryption routines it performs. It also provides mechanisms for decrypting some encrypted values.

The keys used by Anonymizer are randomly generated 128-bit hexadecimal strings. By default, Anonymizer automatically generates a new key on each execution. This is the most secure usage mode as no key is reused more than once. If the generated trace file is to be shared with no plan or need for decryption later, this is a recommended selection.

However, if you – the original trace file owner - should ever want to decrypt encrypted values back into their original form, you will need to retain the key. This can become unwieldy if you use many different keys. To address this, Anonymizer also accepts a pre-defined secret key using the command-line flag `-key` or via the `key` field in the static Options file. This allows for the reuse of keys across many Anonymizer execution runs, ensuring a source IP address is always encrypted to the same value and helps simplify key management.

Finally, to aid in the creation of secret keys, Anonymizer provides a command-line flag `-generate-key` that causes it to simply generate a new key and print it to STDOUT as shown below.

```
./anonymizer -generate-key  
63eba3a49381930fda89b2ad694e74fe
```

While copying and pasting this value is an option, specifying the `-store-options` flag will cause Anonymizer to save a YAML formatted options file containing a dynamically generated key. This is the recommended method if the goal is to produce a static options file, with a static key, for reuse.

Anonymize MAC Addresses in Ethernet Frame Headers

MAC addresses found in the Ethernet header of an ethernet frame uniquely identify the source and destination interfaces of that frame. Depending on where the traffic sample was captured, the source MAC address may belong to the actual transmitting host, or the nearest upstream router from that capturing device. MAC addresses also contain a 3-byte manufacturer prefix that can be used to identify the manufacturer of the network interface card and possibly the device itself. For these reasons, MAC addresses present significant risk of asset identity or information leakage.

The current release of Anonymizer focuses on full anonymization of MAC address values. That is, original MAC address values are replaced with new, randomly generated values that cannot be used to obtain the original value.

If you'd like to disable full MAC anonymization, use the `-mac ignore` flag or set the `mac_addresses` field in the Options file as shown below:

```
data_link:  
  mac_addresses: ignore
```

Remove or Replace VLAN ID Values in 802.1Q Header

While VLAN ID's do not in themselves pose a privacy exposure risk, but they can be used by would-be attackers to gain an understanding of a network's structure and segmentation. For this, VLAN information presents an elevated risk of unintentional information leakage.

Anonymizer offers flexibility in dealing with 802.1Q VLAN headers in the form of a command-line flag `-vlan <mode>`. This command flag allows you to specify one of the following VLAN handling modes:

ignore	Take no action on VLAN headers (default)
remove	Remove VLAN headers, update Ethertype in frame
full	Replace VLAN ID's with pseudorandom number

When **full** mode is specified, Anonymizer randomly selects a new VLAN ID from the range: 2-4093 and replaces the original value with it. This relationship between original and newly selected VLAN ID persists for the remainder of the Anonymizer run.

When **remove** mode is specified, Anonymizer updates the Ethertype of the outer frame header to reflect the type defined in the 802.1Q header, then removes the 802.1Q header entirely. Removing the VLAN header is the recommended mode if a source trace file is known to contain VLAN information that should not be shared. To do this use the `-vlan remove` flag or equivalent Options setting as shown below.

```
data_link:  
  vlan_ids: remove
```

In the case of the Q-in-Q (IEEE 802.1ad) protocol, an Ethernet frame may contain more than one VLAN header or tag. In future versions, Anonymizer will perform recursive anonymization of VLAN tags, however the current recommendation is to omit Q-in-Q frames where possible.

Recalculate Transport and IP Layer Checksums

When modifications are made to an IP header or transport-layer header like UDP or TCP, each header checksum field must be recalculated to reflect the new values. Header values that affect checksum calculation include IP and MAC addresses and since TCP header checksums are calculated using a pseudo-header representing the packet's IP header, changes to the IP header also require the transport header checksum to be recalculated.

Anonymizer's default behavior is to recalculate checksums for all modified packets. However, this action can be disabled using the `-fix-checksums false` flag or equivalent Options setting as shown below:

```
fix_checksums: false
```


Anonymize IP Addresses

IP addresses are among the most important artifacts contained in a PCAP file. They identify the network addresses of devices involved in each recorded communication. IP addresses also pose significant risk of asset and even user identity exposure. For this reason, anonymizing IP addresses effectively is one of Anonymizer's core design goals.

For performing full IP address anonymization, Anonymizer uses *ipcipher*, a specification for encrypting IPv4 and IPv6 addresses. The *ipcipher* specification defines reliable, repeatable mechanisms for encrypting IP addresses, specifically, *ipcrypt* for encrypting IPv4 addresses and AES-128 for encrypting IPv6 addresses.

As *ipcipher* uses deterministic algorithms for encrypting IP addresses, using a static key ensures address anonymization done across multiple Anonymizer runs, against the same source IP addresses, will produce the same encrypted IP address values. However, under most circumstances, using a randomly generated key is the more secure and preferred approach. If you have no intent of decrypting the IP addresses later, the key can be destroyed. Otherwise, the key should be protected and only shared with stake holders with a specific need to possess it.

Anonymizer's IP address anonymization functionality is exposed via the `-ip <mode>` command-line flag and supports the following modes:

full	Perform full anonymization of all IP addresses found in IP headers (default)
ignore	Ignore all IP addresses found in IP headers
include	Anonymize only IP addresses belonging to a subnet defined in <code>-cidr-list</code>
exclude	Anonymize IP addresses that do not belong to a subnet defined in <code>-cidr-list</code>

When **full** mode is specified, Anonymizer encrypts IP addresses found in the IP header of every packet read from the source trace file, this includes IPv4 and IPv6 addresses. This is the default behavior if no IP anonymization mode is specified.

As described in the table above the **include** and **exclude** modes rely on an additional command-line flag `-cidr-list`. The `-cidr-list` flag takes a comma or space-separated list of IP subnets defined in CIDR notation as an argument. Using the IP anonymization mode and CIDR list you can granularly target address anonymization activities to specific IP ranges of interest.

Remove or Anonymize ARP Messages

Address Resolution Protocol (ARP) messages are a fundamental component of Ethernet networks and are used by endpoints and infrastructure devices to resolve an IP address to MAC address and vice versa. ARP messages contain MAC to IP resolution information and can pose a significant information

leakage risk and can completely undermine other IP packet anonymization protections applied to a generated trace file.

For this, we recommend excluding ARP messages from trace files to be shared using the BPF filter expression or by specifying the `-arp remove` flag or equivalent Options setting. However, if you choose to retain ARP messages a full anonymization mode is also available:

full	Perform full anonymization of IP and MAC addresses found in ARP messages
remove	Do not write ARP messages to the generated trace file
ignore	Leave ARP messages in place and unchanged (default)

When **full** mode is selected, Anonymizer uses the same IP and MAC address anonymization functions described above to obfuscate addresses found in an ARP message. Specifically, this includes the Sender and Target Protocol and Hardware address fields. As ARP messages typically contain IP and MAC addresses found in other IP communications, Anonymizer ensures the affected addresses are consistently mapped throughout a generated trace file, regardless of the protocol header.

Remove or Anonymize ICMP Messages

The Internet Control Message Protocol (ICMP) is an IP-layer protocol used for a wide variety of network management and troubleshooting tasks. For example, infrastructure devices such as routers often use ICMP *destination unreachable* or *time exceeded* messages to notify the sender of a packet its intended destination was unreachable. Much like ARP messages, ICMP packets pose significant risk of information leakage as they contain a full IP header as well as a payload. An ICMP payload may contain sensitive data including portions of the sender's original packet.

Like ARP messages, we recommend removing ICMP packets from generated trace files unless they are specifically needed. This can be done using a BPF filter expression, using the `-icmp remove` flag or equivalent Options file setting as shown below.

```
network:
  icmp_messages: ignore
```

If you intend to include ICMP traffic in the generated trace, but wish to obscure the addresses contained with ICMP packets Anonymizer also provides a full ICMP packet anonymization mode:

full	Perform full anonymization of IP addresses found in IP header of ICMP packets
remove	Do not write ICMP packets to the generated trace file
ignore	Leave ICMP packets in place and unchanged (default)

Packet Truncation to User-Specified Length

One of the major limitations of working with raw PCAP files is their size. To help you address the substantial storage requirements associated with working with raw PCAP files, Anonymizer provides the ability to truncate packet payloads that are longer than a specified length.

This option is made available via the `-truncate` flag and the optional `-truncate-len <length>` flag. These options allow you to specify a maximum payload length to which any overly long payload will be truncated. If truncation is enabled but the length is not specified, Anonymizer will truncate overly long packet payloads to a default length of 1024 bytes.

While truncating packet payloads can greatly reduce the size of a generated trace file, doing so does have negative side-effects. Often, 3rd-party analysis tools such as Wireshark will annotate affected packets with messages that describe length mismatches and missing or unseen segments. This is expected behavior as it serves as a signal to the researcher that the packets contained in the trace have been truncated.

BPF Filtering of Packets Read

A very common need – and best practice - when sharing traffic samples is to strip out unwanted traffic. For example, if a researcher has captured all packets on an interface, but they intend to only share HTTP packets, they need an easy way to remove all non-HTTP packets from the trace.

There are many tools that provide this capability when working with raw PCAP files. Anonymizer also accommodates this need by allowing you to specify a Berkley Packet Filter (BPF) expression using the `-bpf <filter expression>` command-line flag or equivalent Options settings as shown below.

```
application:  
  truncate_payload: false  
  truncate_length: 1024
```

The provided filter is compiled and applied when packets are read from the source trace file, before any anonymization tasks are performed. This reduces Anonymizer's overall processing time as well as the size of the generated trace file. We recommend using BPF filters to help ensure only the specific traffic to be shared is included in the generated trace file which greatly helps reduce the risk of accidental information leakage.

Payload Obfuscation

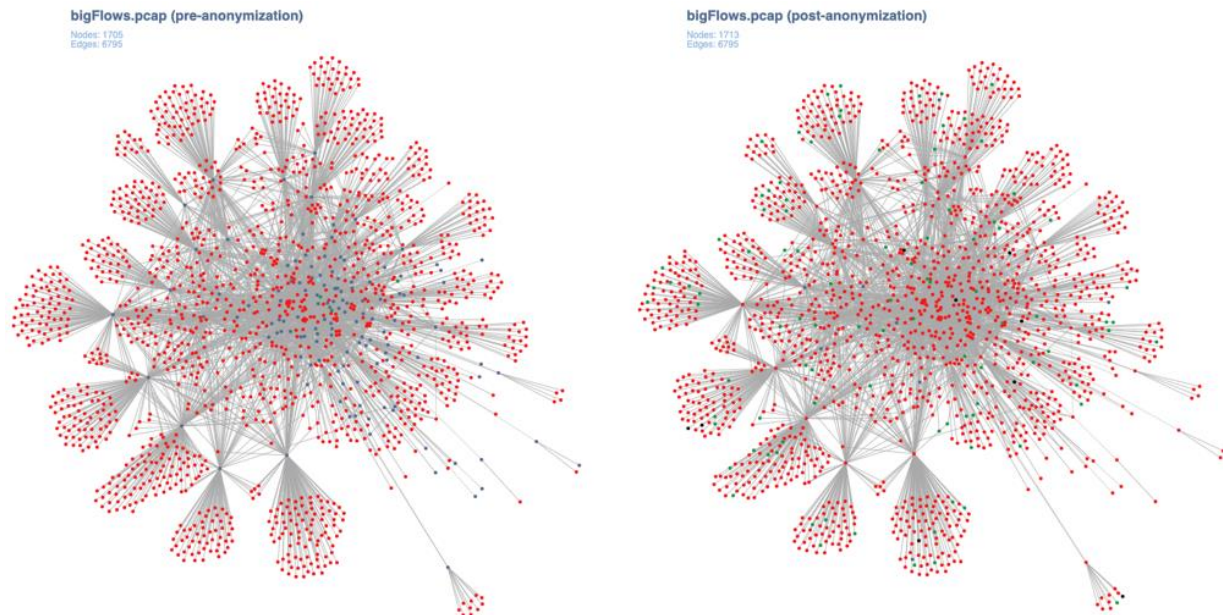
Anonymizer implements a basic payload obfuscation capability whereby payload bytes are replaced with a non-sensical value making the original byte value no longer obtainable. Currently the ASCII asterisk character "*" (hexadecimal 2A) is used to obfuscate the actual payload bytes. To enable this feature, specify the `-mask-payloads` command-line flag or equivalent Options file setting as shown below:

```
application:  
  mask_payload: true
```

Network Graph Views of IP Communications

Through this work we have paid particular attention to maintaining entity-to-entity relationships in the generated trace file. To help validate this, Anonymizer includes a feature that will cause it to generate network graph visualizations of IP communications contained in the trace file, pre and post anonymization. This feature can be enabled by simply passing the `-graph` command-line flag when Anonymizer is executed.

The result is an interactive web page stored in HTML format, in the same directory as the generated trace file. Open the HTML file with a local web browser, and you are presented with a view like that shown below.



The graph on the left represents the IP communications in the trace prior to anonymization. You can zoom in and out, pan the view, and hover over individual nodes to see their IP address value and highlight adjacent edges and connected nodes.

The graph on the right represents IP communications post-anonymization. The same interactive features are available here, allowing you to explore and validate endpoint relationships as well as IP address mappings.

Additional information is encoded in these graphs in the form of glyph shapes and colors. In summary these encodings are:

- Circle Nodes: IPv4 addresses
- Triangle Nodes: IPv6 addresses
- Green Nodes: Multicast and Unicast addresses
- Black Nodes: Loopback addresses
- Light Blue Nodes: IPv4 broadcast addresses
- Red Nodes: IPv4 non-private addresses
- Dark Blue Nodes: IPv4 RFC 1918 (local) addresses