

DetGen: Deterministic Ground Truth Traffic Generation using Docker for Machine Learning

Anonymous Author(s)

ABSTRACT

ACM Reference Format:

Anonymous Author(s). 2018. DetGen: Deterministic Ground Truth Traffic Generation using Docker for Machine Learning. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

The exponential growth of data availability enabled the machine learning revolution of this decade and transformed many areas of our lives. Ironically, researchers struggle to gather qualitative network traffic data to ... Well-designed datasets are such a rarity that researchers often evaluate intrusion detection systems on datasets that are well over a decade old [? ?], calling into question their effectiveness on modern traffic and attacks. The lack of quantity, variability, meaningful labels, and ground truth has so far prohibited ML-based methods from having a bigger impact in network security.

Privacy and security concerns discourage network administrators to release rich and realistic datasets for the public. Network traffic produced by individuals contains a host of sensitive, personal information, such as passwords, email addresses, or usage habits, requiring researchers to expend time anonymising the dataset [?]. In order to examine malicious behaviour, researchers are often forced to build artificial datasets using isolated machines in a laboratory setting to avoid damaging operational devices. Background traffic is generated either from ...

The datasets currently available are meant to be **all-purpose** and are *static* in their design, unable to be modified or expanded. This proves to be a serious defect as the ecosystem of intrusions is continually evolving. Furthermore, it prohibits a more detailed analysis of specific areas of network traffic. To prevent this, new datasets must be periodically built from scratch.

Additionally, **ground truth**

Developing a framework that allows researchers to create datasets that circumvent these issues would be extremely beneficial. We propose that this can be done using Docker [?]. Docker is a service for developing and monitoring containers, also known as OS-level virtual machines. Each Docker container is highly specialised in its purpose, generating traffic related to only a single application process. Therefore, by scripting a variety of Docker-based *scenarios*

that simulate benign or malicious behaviours and collecting the resultant traffic, we can build a dataset with perfect ground truth. Furthermore, these scenarios could be continually enhanced and expanded, allowing for the easy creation of datasets containing modern, up-to-date traffic and attacks.

This is the primary goal of this work. Furthermore, we demonstrate the utility of this framework by performing a series of experiments: one that measures the realism of the network traffic produced by our Docker scenarios and two that would be difficult to perform using a conventional dataset.

2 DATA FORMATS AND EXISTING DATASETS

Computers in a network mostly communicate with each other by sending *network packets* to each other, which are split into the control information, also called packet header, and the user information, called payload. The payload of a packet in general carries the information on behalf of an application and can in be encrypted, while the header contains the necessary information for the correct transmission of the packet, including the transmission protocol layer, IP addresses, etc. Packet-level methods can be divided into payload inspection, header-based, or hybrid. Packets are usually stored in the widespread *pcap* format.

The majority of packets are exchanged between two hosts within bidirectional connections. Another common format of network traffic information is based on connection summaries, also called **network flows**. RFC 3697 [?] defines a network flow as a sequence of packets that share the same source and destination IP address, IP protocol, and for TCP and UDP connections the same source and destination port. A network flow is usually saved containing these informations along with additional information such as the start and duration of the connection as well as the number of packets and bytes transferred in the connection.

3 PROBLEMS IN EXISTING DATASETS AND TRAFFIC GENERATORS

Capturing network traffic into a public dataset is marred by several difficulties and has seen a wealth of criticism. As discussed by Sperotto et al. [?], simply monitoring the usage of several internet users and collecting the resultant traffic into a single dataset, although possible, introduces serious ethical concerns due to the large amount of sensitive or personally identifiable information that average internet users transmit during daily use — such as passwords, GPS coordinates, private material. Despite these difficulties, a few network traffic datasets designed for network intrusion detection exist.

containing a mixture of benign and malicious traffic to test and train machine-learning approaches, captured using tools such as tcpdump or Wireshark.

To avoid this complication, researchers often use statistical approximations of real-world traffic to build such datasets. However,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-9999-9/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

it is unclear to what extent such approximations actually resemble real-world traffic. For instance, it is unclear what the ratio of benign to malicious traffic should be as it is unknown what this ratio is in the real-world. Allix et al. [?] claim that the standard workflow for developing machine-learning systems — namely, collecting large amount of data and then training the algorithm on that data — is necessarily flawed when applied to the domain of intrusion detection. They contend that the inherent secrecy of the intrusion ecosystem and the rate at which it develops make it is impossible to develop a dataset containing, say, network traces of intrusions that are currently being deployed by malicious agents. Instead, one can only build a dataset containing previously discovered attacks. As such, they suggest that it is impossible to release a static dataset that is truly representative of the real-world, impeding the performance of machine-learning-based classifiers.

A problem with all of these datasets is their static design.

Furthermore, many network traffic datasets are not comprehensively labelled. Even if only a single process is initiated by the user, VMs generate additional traffic from a variety of background processes, such as software querying servers to check for updates, resulting in aggregated application flows between processes. This necessitates the use of ground truth generation tools to classify flows. However, current methods of establishing the ground truth of network traffic datasets are well-known to be fallible[?]. These are port-based methods, which classify traffic flows based on their port numbers, and deep-packet inspection (DPI) based methods, which classify flows based on analysis of their packet payloads. Port-based methods are unreliable due to the dynamic allocation of ports, several services sharing the same port and processes running on atypical ports. Moreover, although DPI-based methods are capable of inspecting the entirety of a packets contents, their performance has also shown to be lacking. Bujlow et al. [?] have shown that many commonly used DPI-based methods for ground-truth generation fail to classify packets, with some methods failing to classify common protocols such as HTTP and FTP with less than 10% accuracy. In contrast, it is trivial to produce a fully-labelled dataset from our Docker framework.

CIC-IDS 2017 [?], released by the Canadian Institute for Cybersecurity, is the primary dataset that we shall compare our results to. The dataset was created by monitoring the network activity of several virtual machines running a series of scripted scenarios. The majority of these virtual machines produced exclusively benign traffic whilst others were designated as attackers, producing malicious traffic. Moreover, the exploit scenarios contained within the dataset are moderately recent, including botnets, cross-site scripting attacks and SQL injections. Furthermore, the dataset is far larger than many similar datasets, consisting of five sub-datasets captured over the course of a working week. Capturing traffic over such a lengthy period of time allows for the temporal development of attack scenarios to take place over several days, more accurately mimicking an intruder’s movement through the network. CIC-IDS 2017, however, does not address the problems discussed in this section.

4 REQUIREMENTS

The primary task of this project is to develop a suite of Docker containers capable of producing traffic suitable for training machine-learning-based intrusion detection systems. The containers are arranged in different configurations corresponding to particular *capture scenarios*. Running a given capture scenario triggers the creation of several Docker containers, each with a scripted task specific to that capture scenario. A simple exemplary capture scenario may consist of a containerised client pinging a containerised server. We ensure that each Docker container involved in producing or receiving traffic will be partnered with a tcpdump container, allowing us to collect the resulting network traffic from each container’s perspective automatically. We wish to publish this framework to a wider audience, allowing for further modification. To achieve this goal, we introduce the following key design principles:

- (1) To ensure that we produce representative data for modelling, we want the traffic generated by our container suite to consist of a good number of protocols that are commonly found in real-world traffic and existing datasets. For malicious traffic, we want to ensure that the attacks are modern and varied, both in purpose and in network footprint.
- (2) For each protocol, we want to establish several capture scenarios to encompass the breadth of that protocol’s possible network traces. For instance, if we consider a capture scenario consisting of a client downloading various webpages over SSL, it is not enough to only generate traffic from successful connections. We must also include several scenarios where the client fails to download the aforementioned webpage because of a misspelled address or missing certificate. Whenever possible, we also want to capture WAN traffic.
- (3) Traffic capture scenarios should be implemented in a modular way to allow for a straightforward addition or modification of traffic capture modules.
- (4) Since ground truth a main focus of this work, we want the capture scenarios to be, on some level, deterministic. This way, we can relate individual traffic events to the computational procedure responsible for its generation. We discuss what it means for a scenario to be deterministic in section ??.
- (5) Communication between containers should be subject to the same disturbances and delays as in a real-world setting.

5 VIRTUALISATION AND DOCKER

Virtualisation refers to the act of creating isolated operating systems, known virtual machines (VMs), that share the same hardware infrastructure, known as the host machine. OS-level virtualisation, also known as *containerisation*, is a virtualisation paradigm that has become increasingly popular in recent years due to its lightweight nature and speed of deployment. Containers forego a hypervisor and the shared resources are instead kernel artifacts. Although this prevents the host environment from running different operating systems — for instance, a Linux host can only run Linux containers — containerisation incurs minimal CPU, memory, and networking overhead whilst maintaining a great deal of isolation [?]. The high-level of isolation between both the host OS and any running containers is of paramount importance for our framework as it

allows us to run malicious software that could potentially compromise the security of the system it is running on with negligible worry [?].

Docker container. In Docker’s terminology, a container is a single, running instance of a Docker image. **Add more information about Docker containers, not about images or docker file**

The Docker software platform includes a cloud-based repository called the Docker Hub [?] which allows users to download and build open source images on their local computers. At the time of writing, nearly 2.5 million images are available from Docker Hub.

Docker Networking. Docker allows the creation of virtualised networks with one or more subnetworks, to which containers can securely connect via a virtualised network bridge. Containers attached to the bridge network are assigned an IP address and are able to communicate with other containers on their subnetwork.

Containers can furthermore be connected to a host network, which allows communication with external networks using NAT via the host interface.

We can create our own user-defined Bridge networks, which Docker’s documentation recommends as it provides greater isolation and interoperability between containers [?]. Furthermore, this allows us to set the subnet and gateway for our networks as well as the IP addresses of our containers, which simplifies scripting our scenarios considerably.

Docker-Compose Files. Often, applications built using the Docker framework need more than one container to operate — for example, an Apache server and a MySQL server running in separate containers — and it is therefore necessary to build and deploy several interconnected containers simultaneously. Docker provides this functionality via docker-compose, a tool that allows users to define the services of multiple containers as well as the properties of virtual networks in a YAML file. By default, this file is named *docker-compose.yml*. This allows for numerous containers to be started, stopped and rebuilt with a single command. We can specify volumes, exposed ports, application environment variables, and commands to be run on start up.

6 DESIGN

We provide access to our framework via a GitHub repository, which contains both the implemented capture scenarios as well as the corresponding docker container images.

6.1 Coverage of Protocol Types

As discussed previously, our primary aim is the development of a network traffic dataset for training IDSes. Therefore, in deciding how we should expand the Detlearsom benign scenarios, we initially investigated the protocols and applications present in existing network traffic datasets. Our analysis included CIC-IDS 2017 [?], UNSW-NB15 [?], ICSX Botnet [?] and Mawi [?]. We chose these datasets as they provide .pcap files of their network traffic which enables us to more easily see what protocols are present. To do this, we used the Bro IDS tool to generate log files, listing the results in table ??.

The protocols listed in table ?? make up over 90% of the benign traffic in these datasets. Moreover, although the ratio of protocols

PROTOCOL	UNSW-NB15	ICSX	CIC-IDS 2017	MAWI (2019-7-31)
HTTP	196195	2372	276405	156179
SSL	540	141	285760	591551
DNS	372748	200009	1820105	1581858
X509	459	331	2758590	UNKNOWN
FTP	111685	1989	5540	278
SSH	31320	434	5600	5503
IRC	202	27	0	UNKNOWN
SMTP	44455	125	0	4601

Table 1: Bro Log Flow Count

in datasets can differ significantly, we see some patterns: namely, protocols associated with general browser usage — such as HTTP, SSL, DNS — are the most common in each dataset.

We could base the ratios of protocols in our own dataset off of those found in an existing dataset where the traffic has been artificially generated. However, this would be problematic. In the case of CIC-IDS 2017, some protocols that make up a substantial amount of real-world traffic are glaringly omitted — such as bittorrent or video streaming protocols. In contrast, although UNSW-NB15 contains a better range of protocols, only a small percentage of their web traffic is transmitted over SSL which is completely unrepresentative of real-world traffic. Furthermore, in both these datasets, it appears that a non-negligible amount of traffic is the result of remote desktop protocols being used to interact with the virtual machines generating data, such as X11, VNC and RDP, which we consider to be unwanted noise. Instead, we opt to use these datasets as a rough guideline for constructing datasets but not closely following any one in particular.

6.2 Variation within Protocols

Synthetic datasets often contain traffic from a given protocol, but the manner in which that protocol is used is highly-restricted. It is thus unclear if such traffic is representative of its real-world equivalent usage. For instance, in CIC-IDS 2017 the vast majority of successful FTP transfers consist of a client downloading a .txt file containing the Wikipedia page for ‘Encryption’ several hundred times in a day. This lack of variation in the breadth of a protocol’s network trace results in unrealistic levels of homogeneity in a dataset, and can therefore lead in overoptimistic detection rates of an IDS that do not translate to real-world scenarios. For example, an IDS trained on the CICIDS-17 dataset may trigger a false positive when a larger file is downloaded over FTP due to the difference in the duration or backwards packet length of the connection.

7 CAPTURE

While traversing between to parties, a network packet can pass multiple connecting devices which direct the packet in the right direction. Any device in the immediate circuit traversed by a packet can capture and store it. In a monitoring setting, packets are usually captured by network routers and stored in the widespread *pcap* format. In case of space shortage or privacy concerns, the payload of a packet can be dropped in the saving process.

8 ACKNOWLEDGMENTS

We are grateful for our ongoing collaboration with our industry partners (blinded) on this topic area, who provided both ongoing support and guidance to this work. Discussions with them have

helped reinforce the need for a better evaluation and understanding of the possibilities that new intelligent tools can provide.

Full funding sources after currently blinded.

A RESULT TABLES