

Technical Report CIDDS-002 data set

Markus Ring and Sarah Wunderlich

October 25, 2017

CIDDS-002 (Coburg Intrusion Detection Data Set) [4] is a labelled flow-based port scan data set for evaluation of anomaly based intrusion detection systems. In this report, we provide an overview of the CIDDS-002 data set. We describe in detail the environment in which the data was captured as well as the labelling process of the data set. Further, we explain the structure and the additional published material of the data set. For the underlying ideas of this data set, we refer to our original publication *Creation of Flow-based Data Sets for Intrusion Detection* [4].

Table 1: Revision History

Version	Description	Editors
0.1	First Version of the technical report	Markus Ring, Sarah Wunderlich

1 Terms of Use

To facilitate reproducibility, further development and experiments, we make the CIDDS-002 data set as well as the generation scripts openly available to the community. If you publish material based on our CIDDS (Coburg Intrusion Detection Data Set) dataset or the generation scripts <https://github.com/markusring/CIDDS>, please cite our papers:

Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A.: "Flow-based benchmark data sets for intrusion detection.", In: Proceedings of the 16th European Conference on Cyber Warfare and Security (ECCWS), pp. 361-369. ACPI (2017)

Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A.: "Creation of Flow-based Data Sets for Intrusion Detection.", Information Warfare Journal, vol. 16(4), (2017), to appear

Here is a BiBTeX citation as well:

```
@incollection{ring2017data ,
  title={Flow-based benchmark data sets for intrusion detection},
  author={Ring, Markus and Wunderlich, Sarah and Gruedl, Dominik and
Landes, Dieter and Hotho, Andreas},
  booktitle={Proceedings of the 16th European Conference on Cyber
Warfare and Security (ECCWS)},
  year={2017},
  pages={361--369},
  publisher={ACPI}
}
```

```
@article{ring2017iw ,
  title={Creation of Flow-based Data Sets for Intrusion Detection},
  author={Ring, Markus and Wunderlich, Sarah and Gruedl, Dominik and Landes,
},
  journal={Information Warfare},
  year={2017, to appear},
  volume={16},
  issue={4},
  pages={}
}
```

2 What is the CIDD-002 data set?

CIDD-002 is a labelled flow-based port scan data set for evaluation of anomaly-based network intrusion detection systems. For creation of the CIDD-002 data set, a small business environment was emulated using OpenStack. This environment includes several clients and typical servers like an E-Mail server or a Web server. Python scripts emulate normal user behaviour on the clients.

The CIDD-002 contains unidirectional *NetFlow* [1] data. Table 2 shows an overview of the attributes within the CIDD-001 data set. The attributes 1 to 10 are default *NetFlow* attributes whereas the attributes 11 to 14 are added by us during the labelling process (see Section 5.1).

Table 2: Attributes within the CIDD-002 data set. The second column provides the column names in the published files of the CIDD-002 data set. The third column gives a short description of these attributes.

Nr.	Name	Description
1	Src IP	Source IP Address
2	Src Port	Source Port
3	Dest IP	Destination IP Address
4	Dest Port	Destination Port
5	Proto	Transport Protocol (e.g. ICMP, TCP, or UDP)
6	Date first seen	Start time flow first seen
7	Duration	Duration of the flow
8	Bytes	Number of transmitted bytes
9	Packets	Number of transmitted packets
10	Flags	<i>OR</i> concatenation of all TCP Flags
11	Class	Class label (normal, attacker, victim, suspicious or unknown)
12	AttackType	Type of Attack (portScan, dos, bruteForce, —)
13	AttackID	Unique attack id. All flows which belong to the same attack carry the same attack id.
14	AttackDescription	Provides additional information about the set attack parameters (e.g. the number of attempted password guesses for SSH-Brute-Force attacks)

3 Emulated Network Environment

3.1 Overview

Figure 1 provides an overview of the emulated small business environment in which the CIDDS-002 data set was captured.

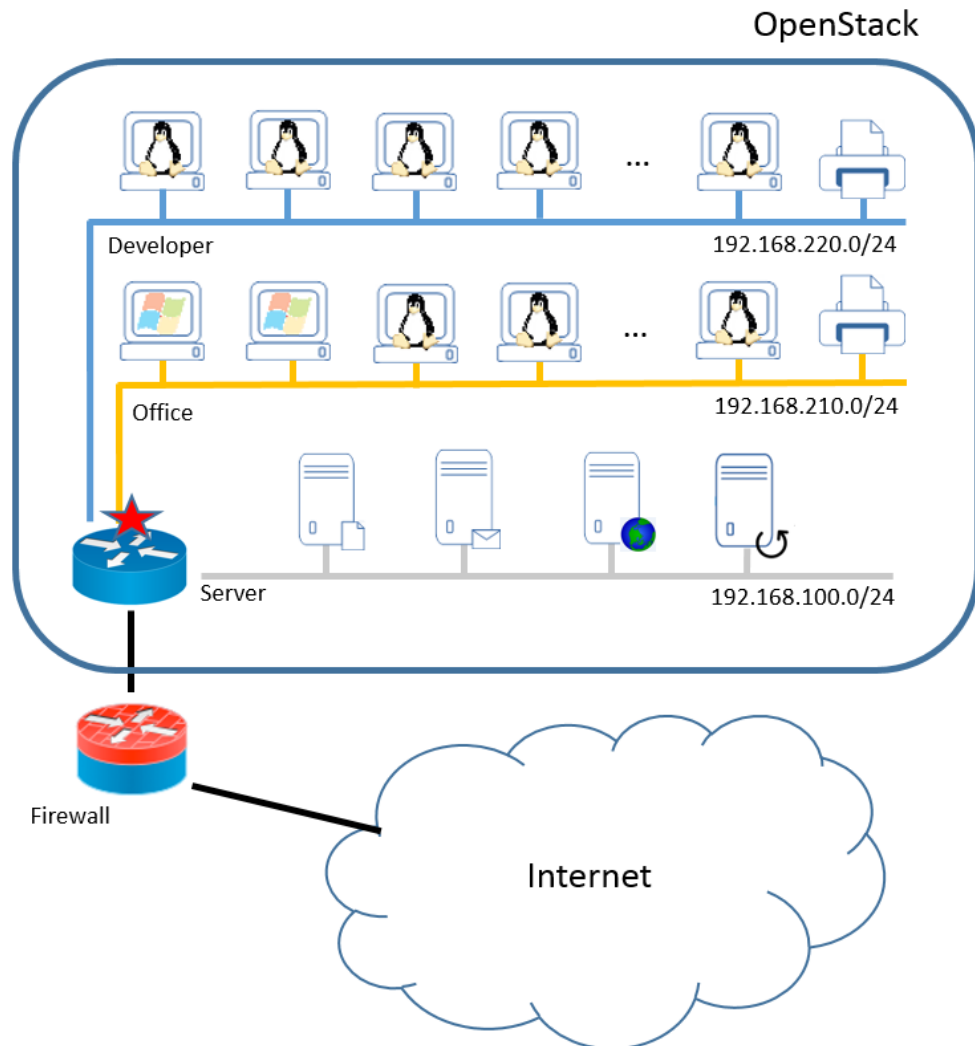


Figure 1: Overview of the small business environment [4]. The star highlights the position where the flow-based traffic was captured.

The unidirectional *NetFlow* traffic was captured at the spot which is highlighted by a red star in Figure 1. The emulated small business environment consists of three subnets. One subnet contains all internal servers (web, file, backup and mail). The other two subnets represent client subnets. Each of them is assigned to a different department (*Developer* and *Office*).

3.2 Server subnet

The server subnet has the Subnet IP *192.168.100.0/24*. It contains four servers (Internal Web server, File server, Mail server and Backup server). The IP Addresses of the servers are shown in Figure 2.

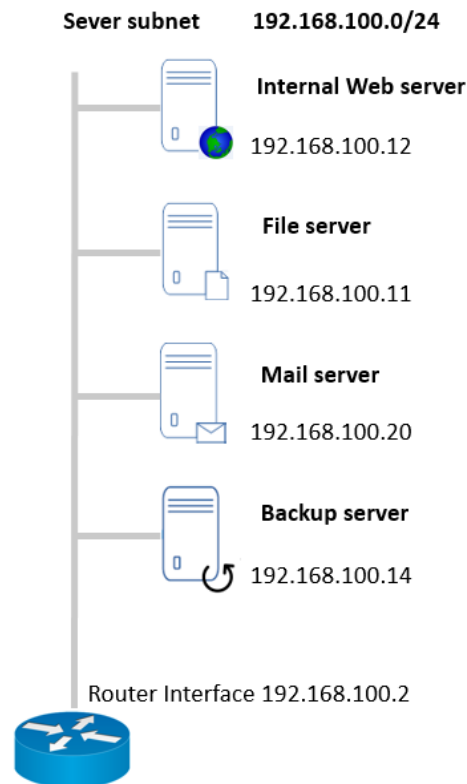


Figure 2: Overview of the *Server* subnet.

The clients occasionally communicate with the *Internal Web server* while surfing the web. From time to time, the clients send and receive E-Mails. In these cases, clients communicate with the *Mail Server*. Further, the clients mount the shared folders from the *File server* as network drives. The backup server is only used by the other three servers. The *Internal Web server*, *Mail server* and *File server* create nightly backups and push them on the *Backup server*.

3.3 Office subnet

The office subnet has the Subnet IP $192.168.210.0/24$. It consists of four *Windows 7* clients, eight *Debian* clients and one network printer. The IP Addresses for the devices are given in Figure 3.

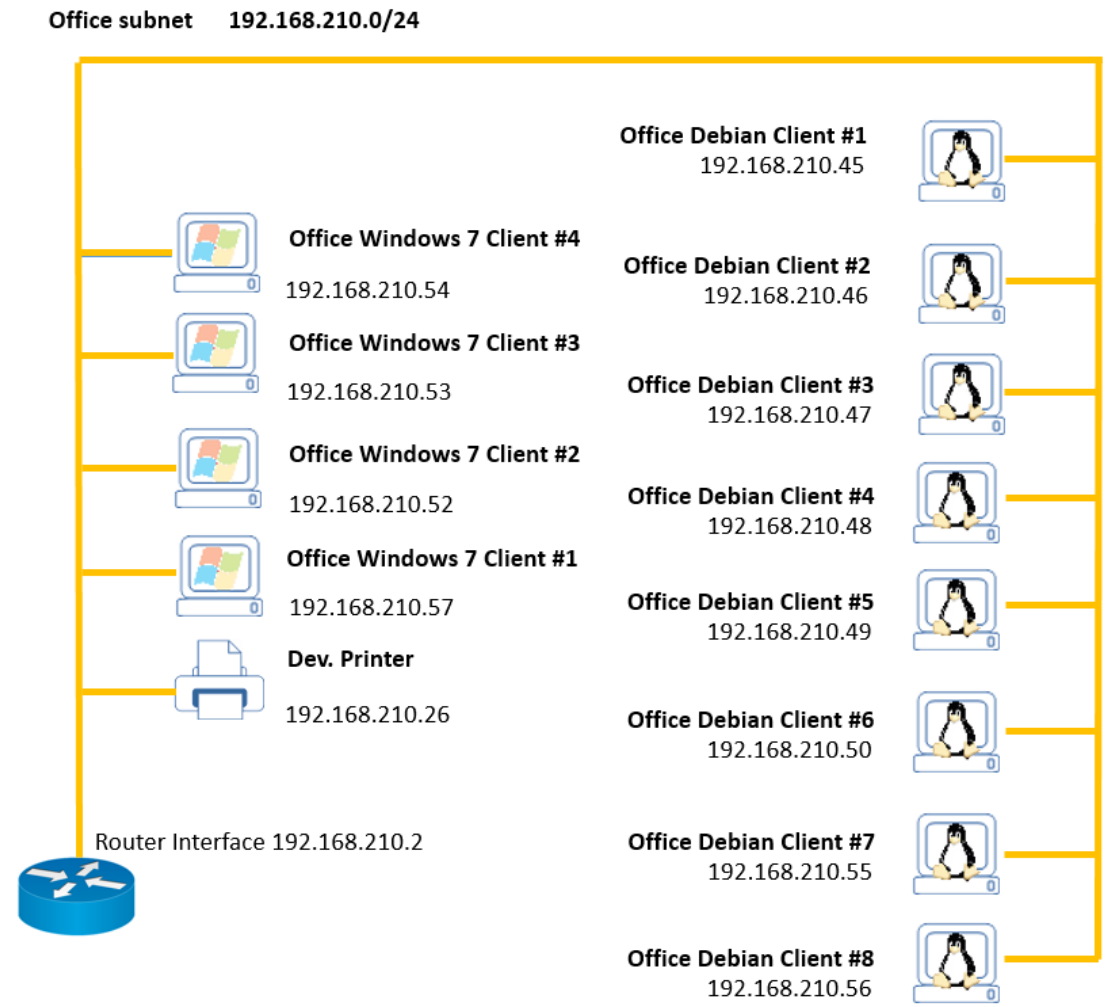


Figure 3: Overview of the *Office* subnet.

3.4 Developer subnet

The *Developer* subnet consists exclusively of *Debian* clients. 10 *Debian* clients and one network printer belong to the subnet IP $192.168.220.0/24$.

The client *Dev. Debian Client #10* executed several port scans within the CIDDS-002 data set. The IP Addresses for the devices are given in Figure 4.

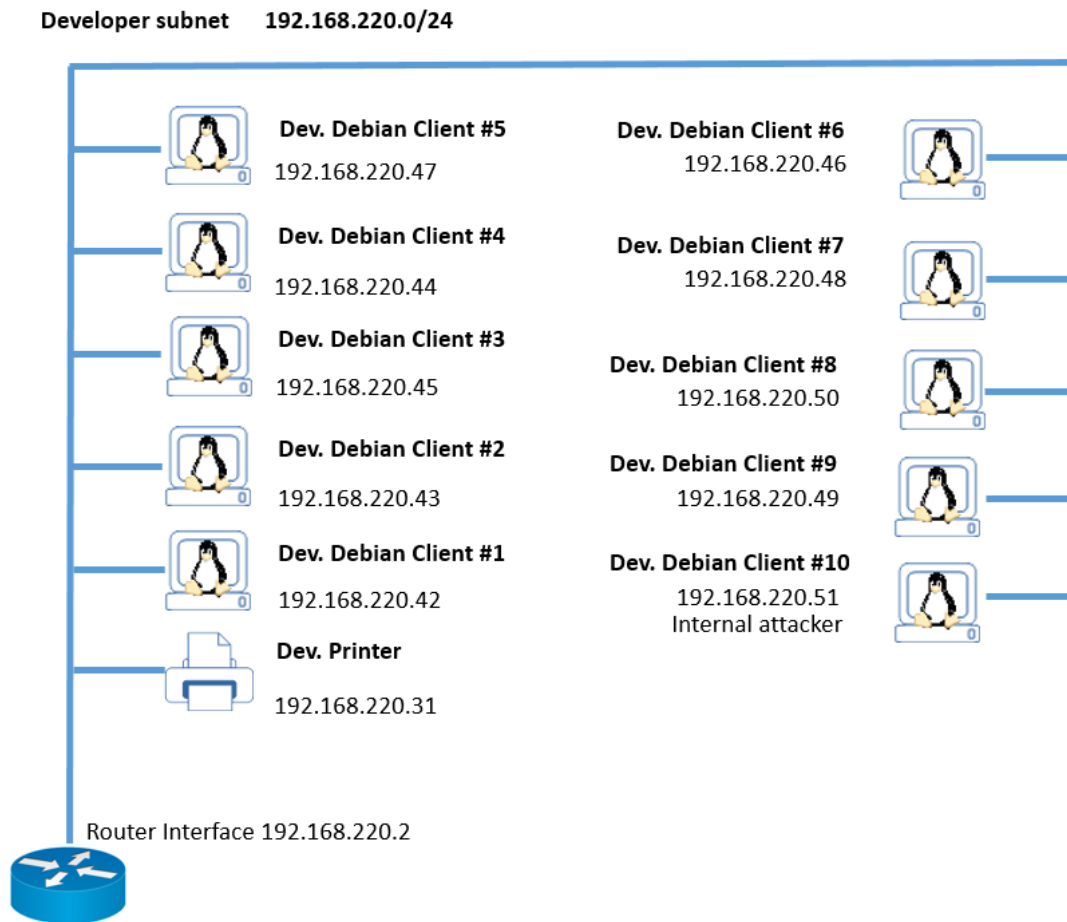


Figure 4: Overview of the *Developer* subnet.

4 Published data of CIDDS-002 data set

This section describes the published material of the CIDDS-002 data set which can be downloaded at <https://www.hs-coburg.de/cidds/>. The archive *CIDDS-002.zip* includes the following four folders: *attack_logs*, *clients_confs*, *client_logs* and *traffic*.

4.1 attack_logs

The folder *attack_logs* contains two *CSV* files which provide information about the executed attacks. Executed attacks within the OpenStack environment are stored in the file *attack_logs_intern.csv*. The file contains information like the *Source IP Addresses*, the *start* and *end time* of the attacks as well as a short description of the attacks.

This file are used for the labelling process in Section 5.1.

4.2 client_confs

As already mentioned above, normal user behaviour of the clients is executed through configurable Python scripts. The *client_confs* folder contains the used configuration files for each client. The file names of the configuration files are constructed as follows:

IP–ADDRESS.conf

The *IP Address* within the file name allows us to map each configuration file to the corresponding client. E.g. the configuration for client *Dev. Debian Client #7* in Figure 4 is *192.168.220.48.conf*.

The configuration file controls – among other things – the different user behaviour of client as well as the working hours of the clients. The user behaviour is defined through the probabilities of the clients activities.

4.3 client_logs

The *client_logs* folder contains for each client a *CSV* file. The name structure of the log files is similar to the structure of the configuration file names:

IP–ADDRESS.log

Remember, the user behaviour of the clients is controlled by Python scripts. These Python scripts record their activities and store them in these log files. This allows users of the *CIDDS-002* data set to understand which client activities caused network traffic.

4.4 traffic

The *traffic* folder contains two *CSV* files with the captured flow-based network traffic in unidirectional *NetFlow* format.

The file names in these sub-folders are constructed as follows:

CIDDS–version–period.csv

All files start with *CIDDS-002*. The last part (*period*) provides information when the network traffic was recorded (*week1* and *week2*).

5 Labelling and Anonymization

We post-processed the recorded data by labelling the flows and by anonymizing public IP Addresses. At first, we describe our labelling process in Section 5.1. Following, the anonymization of all public IP Addresses is described in Section 5.2.

5.1 Labelling

We add four label attributes (*class*, *attackID*, *attackType* and *attackDescription*) to each flow during the labelling process. The first label attribute, called *class*, has three different emphasis: *normal*, *attacker* and *victim*. Each flow is assigned to a predefined class.

The other three label attributes provide additional information about executed attacks. These attributes are only used if the flow belongs to the class *attacker* or *victim*. If the flow belongs to the class *normal*, *suspicious* or *unknown*, the value of the additional label attributes are set to a default value (" - -"). The second label attribute is called *attackID*. In the CIDD-002 data set, a unique ID is assigned to each executed attack. Consequently, all flows which belong to the same attack share the same value in this attribute.

The third label attribute is called *attackType* and gives more information about the executed attack type. Possible values are: *pingScan*, *portScan*, *bruteForce* or *dos*.

The fourth label attribute contains detailed information about the executed attacks. For example, this label attribute contains the parameters settings for *portScans* or the number of passwords guessed for *bruteForce* attacks.

We have full control over the virtual machines and virtual networks within the *OpenStack* environment. This allows us to label all flows with their corresponding classes. Since we know the exact timestamp, origin and target of executed attacks, we are able to label all flows which are caused by attacks. All remaining flows within the *OpenStack* environment are labelled as *normal*.

5.2 Anonymization

We anonymized all public IP Addresses for privacy reasons. The IP Addresses of the internal OpenStack clients and servers are not affected during the anonymization process.

5.2.1 Special Treatments

The following IP Addresses are handled specifically during the anonymization process:

- All virtual machines within the OpenStack environment use the same DNS server. We rename the IP Address of the DNS server as *DNS*.

5.2.2 Other IP Addresses

We used the following anonymization process for the remaining public IP Addresses. The first three bytes of each IP Address are replaced with a randomly generated number. The fourth byte of the IP Address is kept. This allows us to retain information about network structures, since all IP Addresses from the same subnet are replaced with the same randomly generated number. Table 3 shows a few examples for this anonymization process.

6 Traffic Characteristics

The CIDD-001 data set was captured over a period of two weeks and contains about 16 millions flows.

Table 3: Anonymization process of public IP Addresses.

#	IP Address	Anonymized IP Address
1	8.8.8.8	4711_8
2	8.8.8.9	4711_9
3	8.8.8.18	4711_18
4	9.9.9.9	13107_9
5	9.9.9.173	13107_173
6	8.8.8.9	4711_9
7	7.7.7.7	2311_7

The CIDDS-001 data set includes 43 attacks. 20 attacks were executed in *week1* and 13 attacks were executed in *week2*. Table 4 provides more information about the executed attacks within CIDDS-002 data set.

Table 4: Overview of the number and type of executed port scans within the CIDDS-002 data set. Each row represents the attacks within a specific week. The columns describe the different types of port scans.

week	SYN Scan -T1	SYN Scan -T2	SYN Scan -T3	ACK Scan -T1	ACK Scan -T2	ACK Scan -T3	UDP Scan -T0	UDP Scan -T1	UDP Scan -T2	UDP Scan -T3	FIN Scan -T1	FIN Scan -T2	FIN Scan -T3	Ping Scan -T1	Ping Scan -T2	sum
week1	2	1	1	2	1	2	1	2	2	0	0	2	3	0	1	20
week2	3	2	0	1	0	1	0	1	3	1	3	3	2	2	1	23
overall	5	3	1	3	1	3	1	3	5	1	3	5	6	2	2	43

As can be seen in Table 4, week one and week two contain traffic with benign behaviour as well as attacks.

There is one peculiarity to be noted when analyzing the CIDDS-002 data set. In *week1* there is missing traffic between 03:52 AM to 09:08 AM on August 04th due to a failure in OpenStack.

7 Further Information

For further information, we encourage you to read our paper *Creation of Flow-based Data Sets for Intrusion Detection* [2]. In addition to that, we also published our generation and labelling scripts in a github repository [3].

Acknowledgements

This work is funded by the Bavarian Ministry for Economic affairs through the WISENT project (grant no. IUK 452/002).

References

- [1] Claise, B.: Cisco systems netflow services export version 9. RFC 3954 (2004)
- [2] Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A.: Flow-based benchmark data sets for intrusion detection. In: Proceedings of the 16th European Conference on Cyber Warfare and Security (ECCWS), to appear. ACPI (2017)
- [3] Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A.: Generation scripts for the coburg intrusion detection data sets (cidds) (Apr 2017), <https://github.com/markusring/CIDDS>
- [4] Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A.: Creation of flow-based data sets for intrusion detection. Information Warfare 16 (2017, to appear)