

Dataset: Large-scale Urban IoT Activity Data for DDoS Attack Emulation

Arvin Hekmati
hekmati@usc.edu
University of Southern California
Los Angeles, California, USA

Eugenio Grippo
egrippo@usc.edu
University of Southern California
Los Angeles, California, USA

Bhaskar Krishnamachari
bkrishna@usc.edu
University of Southern California
Los Angeles, California, USA

ABSTRACT

As IoT deployments grow in scale for applications such as smart cities, they face increasing cyber-security threats. In particular, as evidenced by the famous Mirai incident and other ongoing threats, large-scale IoT device networks are particularly susceptible to being hijacked and used as botnets to launch distributed denial of service (DDoS) attacks. Real large-scale datasets are needed to train and evaluate the use of machine learning algorithms such as deep neural networks to detect and defend against such DDoS attacks. We present a dataset from an urban IoT deployment of 4060 nodes describing their spatio-temporal activity under benign conditions. We also provide a synthetic DDoS attack generator that injects attack activity into the dataset based on tunable parameters such as number of nodes attacked and duration of attack. We discuss some of the features of the dataset. We also demonstrate the utility of the dataset as well as our synthetic DDoS attack generator by using them for the training and evaluation of a simple multi-label feed-forward neural network that aims to identify which nodes are under attack and when.

CCS CONCEPTS

• Security and privacy → Malware and its mitigation.

KEYWORDS

DDoS Attacks, datasets, neural networks, machine learning, botnet

ACM Reference Format:

Arvin Hekmati, Eugenio Grippo, and Bhaskar Krishnamachari. 2018. Dataset: Large-scale Urban IoT Activity Data for DDoS Attack Emulation. In *SenSys '21: ACM Conference on Embedded Networked Sensor Systems, November 03–05, 2021, SenSys, Coimbra, Portugal*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Technological evolution has made possible the deployment of large internet of thing (IoT) systems that are able to connect multiple different sensors and actuators, allowing them to communicate and

exchange enormous amounts of data. Such large-scale IoT systems consisting of thousands of sensor nodes are being proposed, for example, in the context of smart-city applications ([1], [2], [3]). As these IoT systems grow in size and complexity, they are increasingly vulnerable from a Cybersecurity perspective ([4], [5], [6]).

Most significantly, they are liable to be hacked into and hijacked by malicious entities and then used as part of massive botnets as a launching ground for distributed denial of service (DDoS) attacks, potentially affecting millions of end-users ([7], [8], [9], [10], [11], [12]). A famous example of an IoT-based DDoS attack was the Mirai botnet, first identified in August 2016 by MalwareMustDie, a whitehat security research group. Afterward, some of the biggest DDoS attacks in history were performed by Mirai botnet and its mutated variants. 400,000 nodes infected by this malware executed DDoS attacks on websites with a massive peak of 1.1 Tbps data transfer ([13], [14], [15], [16], [17], [18], [19]).

One of the major steps to defend against such DDoS attacks is to detect them successfully, as close to real-time as possible. To this end, security researchers have turned to machine learning tools such as deep learning networks for DDoS attack detection ([20]). The latest development in technology (parallel computing, high computational processing speeds, GPU's, etc.) have made possible the vertiginous development of deep neural networks (NN's) ([21], [22], [23]), however their performance is very much dependent on the availability of rich datasets to train them. The training of NN models will be challenged in the near future as botnet attacks become increasingly complex and grow in volume infecting massive portions of IoT networks. Therefore, it is important to build and maintain suitable large-scale IoT dataset repositories, useful for training such models, to allow the research community to stay ahead of malicious attack developments.

There have been some prior efforts on collecting or creating synthetic DDoS attack datasets that could be used for training data-driven deep learning models (e.g., [24], [25], [26], [27], [28], [29]). However many of these studies either are not specifically focused on IoT devices (e.g., [30], [31], [32], [33], [34], [35], [36], [37]), while others present data from a limited number of nodes (e.g., [38], [39], [40]). Further, while today's DDoS attacks such as Mirai are often relatively easy to detect because their traffic volumes dramatically exceed normal traffic (e.g., [41] reports nearly 100% true positives for an autoencoder-based detection scheme on the Mirai botnet), we believe that future attacks will be more cleverly camouflaged so that the traffic generated during the attack from a given node matches the traffic volume from a benign node. Thus, beyond datasets that show the behavior of today's attacks, there is still a need for large-scale datasets that characterize the benign activity of real IoT devices, which could be used as a basis to emulate more challenging future attacks.

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.

SenSys '21, November 03–05, 2021, SenSys, Coimbra, Portugal

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/10.1145/1122445.1122456>

We present in this work first a data set obtained from a real urban IoT system in a large city consisting of more than 4000 spatially distributed sensors. The data consists of the binary activity status of each node at a granularity of 30 seconds over a period of one month under a benign (non-attacked) setting.

To make the dataset useful for training machine learning tools for DDoS detection, we need to augment it through synthetic attack emulation. Therefore, in addition to the raw (benign) activity data, we also provide a dedicated script that generates attacks in the proposed dataset synthetically. Our script allows the setting of multiple parameters: number of nodes to be processed, total attack duration, attack ratio, starting time of the attack; and each particular node is equipped with a time-stamp and an output that varies binary between 0 (no attack) and 1 (attacked node). To illustrate the utility of our dataset and attack emulator, we design, train, and implement a simple supervised feed-forward NN model to detect malicious attacks utilizing the provided dataset. We make the dataset and the attack emulation script along with our illustrative NN model available as an open-source repository online at https://github.com/ANRGUSC/Urban_IoT_DDoS_Data.

The paper is structured as follows: section 2 presents the original dataset and some statistics about it. The attack and defense mechanism are presented in section 3. Finally, we conclude the paper in section 4.

2 ORIGINAL AND BENIGN ACTIVITY DATASETS

The original data has been collected from the activity status of real event-driven IoT nodes deployed in an urban area¹. The original dataset contains three main features, the node ID, the location of the node in Latitude and Longitude, and a timestamp of the activity status of the IoT node. A record has been added to the original dataset whenever the activity status of a node changes. The raw dataset has 4060 nodes with one month worth of data.

Having a record of each node whenever the status of that node changes provides a bias towards the information of nodes that have more activity changes during the day. In order to overcome this issue, we also provide a script that takes the original dataset and generates a new benign activity dataset showing the activity status for each node every t_s seconds. In this way, all nodes in the benign activity dataset will have the same number of records. The script can generate a customized benign dataset by providing the beginning and ending date, the number of IoT nodes, and the time step, i.e., t_s .

2.1 Dataset Statistics

In this subsection, some statistics of the dataset are presented to illustrate the dataset's properties.

Figure 1 presents the mean number of active nodes versus the time of the day on one particular day of the dataset. As we can see, up to 65% of the nodes get activated around the middle of the day, but by midnight only about 20% of the nodes are active.

Figure 2 shows the mean correlation between pairs of 500 randomly selected nodes in the dataset, in terms of their activity versus

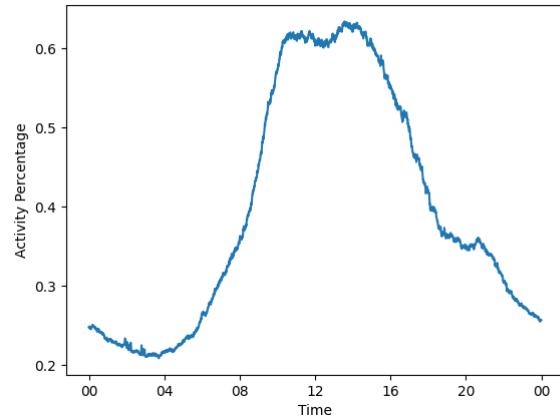


Figure 1: Active Nodes Percentage vs Time

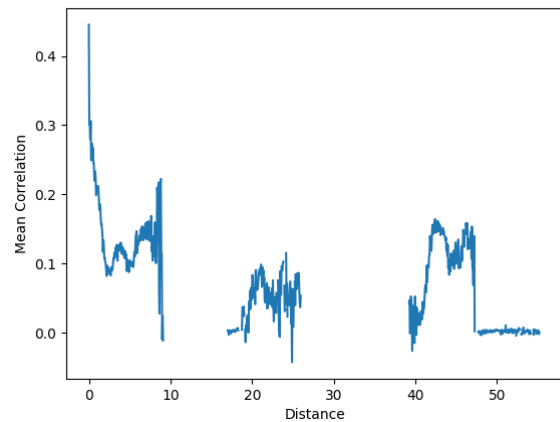
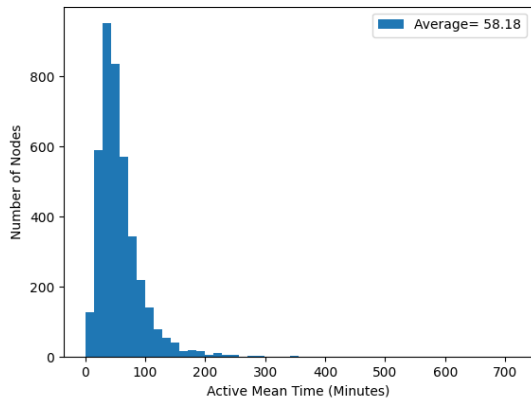


Figure 2: Nodes Activity Mean Correlation vs Distances

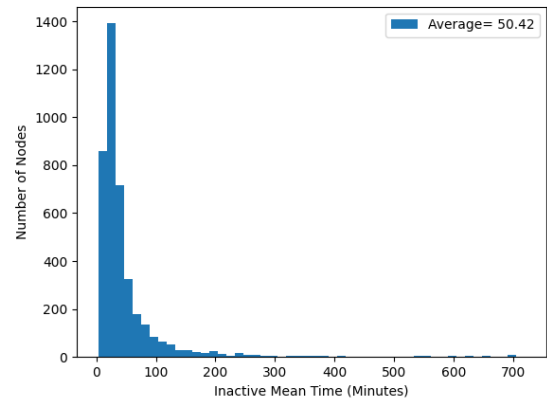
their distance. For this purpose, the pair distances between the nodes are calculated using Euclidean distance. The Pearson correlation has been used to calculate the correlation between nodes' activity in the one month worth of data. Then, the distances between nodes have been split into 1000 bins, and the average correlation for the nodes whose pair distances fall in the bin boundaries has been calculated for each bin. As we can see, generally, the farther the nodes are, their correlation will be lower (dropping from a mean correlation of about 0.4 for nodes that are very close to each other to values around 0.1 or below for distances greater than 10 units).

Next, we present some statistics that show how long nodes tend to remain active or inactive. Figure 3 shows a histogram of the mean active time for all nodes as well as the mean inactive time, for both the day time (8 AM to 8 PM) and nighttime (8 PM to 8 AM). Figures 3a and 3c show the active mean time for the day and night, respectively. Figures 3b and 3d show the inactive mean time for the day and night, respectively. As we can see in these figures, both

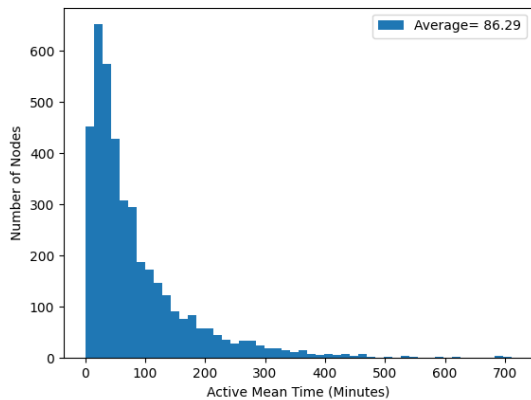
¹The source of this data has been anonymized for privacy and security reasons.



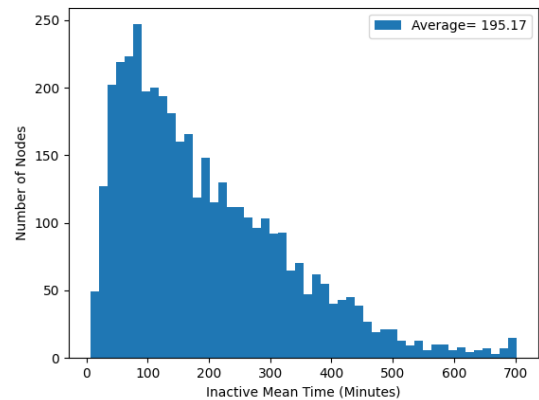
(a) Histogram of mean activity time per node - From 8 AM to 8 PM



(b) Histogram of mean inactivity time per node - From 8 AM to 8 PM



(c) Histogram of mean activity time per node - From 8 PM to 8 AM



(d) Histogram of mean inactivity time per node - From 8 PM to 8 AM

Figure 3: Histograms of mean activity and inactivity time per node

mean active times and mean inactive times tend to be higher in the nights compared to the days.

3 ATTACK AND DEFENSE MECHANISM

This section presents how synthetic DDoS attacks are generated on the IoT nodes. Furthermore, here we define the training dataset features and also the detection mechanism.

3.1 Generating Attack Dataset

In this paper, we synthetically generate a DDoS attack on the IoT nodes by setting all attacked nodes to an active status for the duration of the attack. Note that this approach to attack generation is more coarse-grained than providing the volume of packets or packets generated with specific content and destinations during an attack. We plan to incorporate the generation of such additional fine-grained information during the attacks in the near-future. But

even by focusing on activity status only, we are effectively emulating more challenging futuristic DDoS attacks that may be hard to detect from a single node's traffic.

A script is provided that can be used for generating attacks over the dataset. Three parameters can be set in generating the attacks: the start time of the attack, the duration of the attack, and the percentage of the nodes that go under attack. In our experiments, we used one week's worth of dataset for generating attacks. The attacks are started at 2 AM on each day of the week over all of the nodes with durations of 1, 2, 4, 8, 16 hours.

3.2 Generating Training Dataset

Given the attacked dataset, a labeled training dataset will be generated by calculating the mean activity time of the IoT nodes in the specified time windows. Note that, in the script provided for this paper, one could set the desired time windows for generating the training features. In this paper, we considered a list of 12 different

Table 1: Mean accuracy and recall for the NN detection model

| | Mean Accuracy | Mean Recall |
|------------------|---------------|-------------|
| Training Dataset | 0.94 | 0.93 |
| Testing Dataset | 0.88 | 0.84 |

time windows, namely as 1, 10, 30, minutes and 1, 2, 4, 8, 16, 24, 30, 36, 42, 48 hours, to calculate the mean activity time of the nodes.

3.3 Defense Mechanism

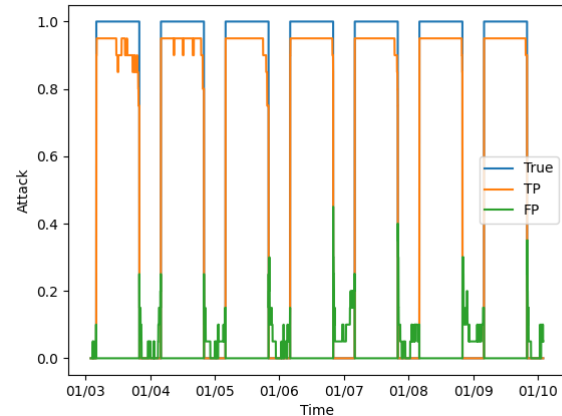
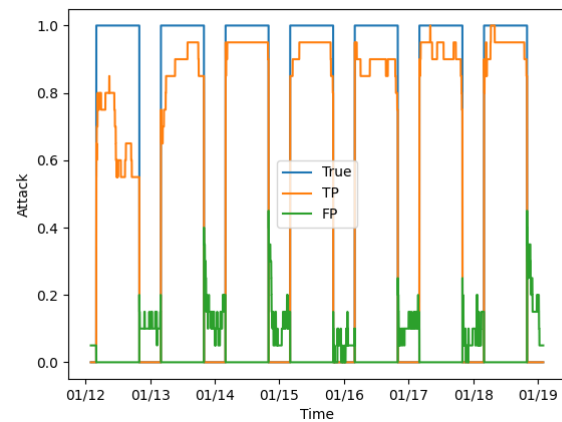
We train a feed-forward neural network to detect the DDoS attack on the IoT nodes based on their collective activity status over time. As noted before, detecting attacks based on activity status alone is more challenging than approaches based on measuring fine-grained traffic volumes or flow-level information.

The model we have trained is a simple binary classification to show a sample usage of the presented dataset. In this neural network, we have an input layer with 12 neurons. The input layer is followed by one hidden layer with 8 neurons and ReLU activation. A dropout of 20% and batch normalization is also used at the end of the hidden layer. The output is a single neuron with the Sigmoid activation function. In this experiment, we randomly selected 20 IoT devices nodes to generate attacks. We are training 20 different models for each IoT node using its data alone, each with different weights but all having the same architecture. The neural network model is trained for 500 epochs for each node to detect the attacked time slots in the dataset. We used one week's worth of data as the training dataset and another week as the testing dataset. Note that this is a simple approach that will not take into account any correlations in the data across different nodes, so there is scope for further improvement by developing more complex models that integrate the inputs from multiple IoT devices.

Table 1 present the mean recall and accuracy of the 20 models trained for detecting DDoS attacks. Figures 4 and 5 show the true attack attack (T), attack predictions true positive (TP) and false positives (FP) mean over all nodes vs time, for both training and testing dataset. In these figures we used the attack duration of 16 hours. As we can see, the attacked nodes are being detected very well in the training dataset with a few FP. On the other hand, we are getting around 84% recall on the testing dataset with a little bit more FP.

4 CONCLUSION

We have presented a new spatio-temporal dataset describing the activity of a 4060-node event-based urban IoT deployment. We have also provided a script to create a benign dataset out of the original dataset to reduce bias toward nodes with more activity. We have shown some statistical analyses on the dataset to illustrate some its key properties. We have also presented a synthetic DDoS attack generator to generate attacks using the dataset, and illustrate the training and evaluation of a feed-forward neural network using the dataset as a way to detect such attacks. We hope that the dataset and tools we have provided are helpful to the community to undertake various types of research related to large-scale event-driven IoTs, including DDoS attack detection.

**Figure 4: Training Dataset Attack Prediction vs Time****Figure 5: Testing Dataset Attack Prediction vs Time**

5 ACKNOWLEDGMENTS

This material is based upon work supported by Defense Advanced Research Projects Agency (DARPA) under Contract No. HR001120C0160 for the Open, Programmable, Secure 5G (OPS-5G) program. Any views, opinions, and/or findings expressed are those of the author(s) and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

REFERENCES

- [1] H. Arasteh, V. Hosseinezhad, V. Loia, A. Tommasetti, O. Troisi, M. Shafie-khah, and P. Siano. Iot-based smart cities: A survey. In *2016 IEEE 16th International Conference on Environment and Electrical Engineering (EEEIC)*, pages 1–6, 2016.
- [2] Camilo Alejandro Medina, Manuel Ricardo Pérez, and Luis Carlos Trujillo. Iot paradigm into the smart city vision: A survey. In *2017 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, pages 695–704, 2017.
- [3] Ahmed Samy Nassar, Ahmed Hossam Montasser, and Nashwa Abdelbaki. A survey on smart cities' iot. In *International Conference on Advanced Intelligent Systems and Informatics*, pages 855–864. Springer, 2017.

- [4] Yang Lu and Li Da Xu. Internet of things (iot) cybersecurity research: A review of current research topics. *IEEE Internet of Things Journal*, 6(2):2103–2115, 2018.
- [5] Dhuha Khalid Alferidah and NZ Jhanjhi. Cybersecurity impact over bigdata and iot growth. In *2020 International Conference on Computational Intelligence (ICCI)*, pages 103–108, 2020.
- [6] Ibrahim Alrashdi, Ali Alqazzaz, Esam Aloufi, Raed Alharthi, Mohamed Zohdy, and Hua Ming. Ad-iot: Anomaly detection of iot cyberattacks in smart city using machine learning. In *2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 0305–0310, 2019.
- [7] Elisa Bertino and Nayeem Islam. Botnets and internet of things security. *Computer*, 50(2):76–79, 2017.
- [8] Constantinos Kolias, Georgios Kambourakis, Angelos Stavrou, and Jeffrey Voas. Ddos in the iot: Mirai and other botnets. *Computer*, 50(7):80–84, 2017.
- [9] Roger Hallman, Josiah Bryan, Geancarlo Palavicini, Joseph Divita, and Jose Romero-Mariona. Iodds-the internet of distributed denial of service attacks. In *2nd international conference on internet of things, big data and security. SCITEPRESS*, pages 47–58, 2017.
- [10] Bruno Bogaz Zarpelo, Rodrigo Sanches Miani, Cludio Toshio Kawakani, and Sean Carliso de Alvarenga. A survey of intrusion detection in internet of things. *J. Netw. Comput. Appl.*, 84(C):25–37, April 2017.
- [11] Sebastián García, Alejandro Zunino, and Marcelo Campo. Survey on network-based botnet detection methods. *Security and Communication Networks*, 7(5):878–903, 2014.
- [12] Saman Taghavi Zargar, James Joshi, and David Tipper. A survey of defense mechanisms against distributed denial of service (ddos) flooding attacks. *IEEE communications surveys & tutorials*, 15(4):2046–2069, 2013.
- [13] Manos Antonakakis, Tim April, Michael Bailey, Matt Bernhard, Elie Bursztein, Jaime Cochran, Zakir Durumeric, J Alex Halderman, Luca Invernizzi, Michalis Kallitsis, et al. Understanding the mirai botnet. In *26th {USENIX} security symposium ({USENIX} Security 17)*, pages 1093–1110, 2017.
- [14] Hamdija Sinanović and Sasa Mrdovic. Analysis of mirai malicious software. In *2017 25th International Conference on Software, Telecommunications and Computer Networks (SoftCOM)*, pages 1–5. IEEE, 2017.
- [15] K Vengatesan, Abhishek Kumar, M Parthibhan, Achintya Singhal, and R Rajesh. Analysis of mirai botnet malware issues and its prediction methods in internet of things. In *International conference on Computer Networks, Big data and IoT*, pages 120–126. Springer, 2018.
- [16] Constantinos Kolias, Georgios Kambourakis, Angelos Stavrou, and Jeffrey Voas. Ddos in the iot: Mirai and other botnets. *Computer*, 50(7):80–84, 2017.
- [17] Christopher D. McDermott, Farzan Majdani, and Andrei V. Petrovski. Botnet detection in the internet of things using deep learning approaches. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2018.
- [18] Joel Margolis, Tae Tom Oh, Suyash Jadhav, Young Ho Kim, and Jeong Neyo Kim. An in-depth analysis of the mirai botnet. In *2017 International Conference on Software Security and Assurance (ICSSA)*, pages 6–12, 2017.
- [19] Artur Marzano, David Alexander, Osvaldo Fonseca, Elverton Fazzion, Cristine Hoepers, Klaus Steding-Jessen, Marcelo H. P. C. Chaves, Ítalo Cunha, Dorgival Guedes, and Wagner Meira. The evolution of bashlite and mirai iot botnets. In *2018 IEEE Symposium on Computers and Communications (ISCC)*, pages 00813–00818, 2018.
- [20] Niharika Sharma, Amit Mahajan, and V Malhotra. Machine learning techniques used in detection of dos attacks: a literature review. *International Journal of Advance Research in Computer Science and Software Engineering*, 6(3):100–105, 2016.
- [21] Kevin Gurney. *An introduction to neural networks*. CRC press, 2018.
- [22] Ethem Alpaydin. *Introduction to machine learning*, 3rd editio. ed, 2014.
- [23] Amer Zayegh and N Bassam. *Neural network principles and applications*, volume 10. London: Pearson, 2018.
- [24] University of New Brunswick, Canadian Institute for Cybersecurity. <https://www.unb.ca/cic/datasets/>. Accessed: 09-11-2021.
- [25] University of California Irvine, KDD Archive. <https://www.kdd.org/kdd-cup/view/kdd-cup-1999/Data>. Accessed: 09-11-2021.
- [26] University of New South Wales, The UNSW-NB15 Dataset. <https://research.unsw.edu.au/projects/unsw-nb15-dataset>. Accessed: 09-11-2021.
- [27] Dilara Gümüşbaş, Tulay Yıldırım, Angelo Genovese, and Fabio Scotti. A comprehensive survey of databases and deep learning methods for cybersecurity and intrusion detection systems. *IEEE Systems Journal*, 15(2):1717–1731, 2021.
- [28] Markus Ring, Sarah Wunderlich, Dominik Gründl, Dieter Landes, and Andreas Hotho. Flow-based benchmark data sets for intrusion detection. In *Proceedings of the 16th European Conference on Cyber Warfare and Security. ACPI*, pages 361–369, 2017.
- [29] Markus Ring, Sarah Wunderlich, Dominik Gründl, Dieter Landes, and Andreas Hotho. Creation of flow-based data sets for intrusion detection. *Journal of Information Warfare*, 16:40–53, 2017.
- [30] The CAIDA UCSD "DDoS Attack 2007" Dataset. https://www.caida.org/catalog/datasets/ddos-20070804_dataset. Accessed: 09-12-2021.
- [31] Iman Sharafaldin, Arash Habibi Lashkari, Saqib Hakak, and Ali A Ghorbani. Developing realistic distributed denial of service (ddos) attack dataset and taxonomy. In *2019 International Carnahan Conference on Security Technology (ICCSST)*, pages 1–8. IEEE, 2019.
- [32] DARPA intrusion detection evaluation dataset, Lincoln Laboratory, Massachusetts Institute of Technology. <https://www.ll.mit.edu/r-d/datasets/2000-darpa-intrusion-detection-scenario-specific-datasets>. Accessed: 09-12-2021.
- [33] University of New Brunswick, "DDoS Evaluation Dataset (CICDDoS2019)", unb.ca, 2019. <https://www.unb.ca/cic/datasets/ddos-2019.html>. Accessed: 09-13-2021.
- [34] University of New Brunswick, "CSE-CIC-IDS2018 on AWS", 2018. <https://www.unb.ca/cic/datasets/ids-2018.html>. Accessed: 09-13-2021.
- [35] Canadian Institute for Cybersecurity, "CICIDS2017", unb.ca, 2017. <https://www.unb.ca/cic/datasets/ids-2017.html>. Accessed: 09-13-2021.
- [36] Frank Beer, Tim Hofer, David Karimi, and Ulrich Bühler. A new attack composition for network security. In Paul Müller, Bernhard Neumair, Helmut Raiser, and Gabi Dreo Rodosek, editors, *10. DFN-Forum Kommunikationstechnologien*, pages 11–20, Bonn, 2017. Gesellschaft für Informatik e.V.
- [37] Derya Erhan and Emin Anarım. Boğaziçi university distributed denial of service dataset. *Data in brief*, 32:106187, 2020.
- [38] A Scheme for Generating a Dataset for Anomalous Activity Detection in IoT Networks. <https://sites.google.com/view/iot-network-intrusion-dataset>. Accessed: 09-12-2021.
- [39] University of California Irvine, Machine Learning Repository. https://archive.ics.uci.edu/ml/datasets/detection_of_IoT_botnet_attacks_N_BaIoT. Accessed: 09-12-2021.
- [40] University of New South Wales, The Bot-IoT Dataset. <https://research.unsw.edu.au/projects/bot-iot-dataset>. Accessed: 09-12-2021.
- [41] Yair Meidan, Michael Bohadana, Yael Mathov, Yisroel Mirsky, Asaf Shabtai, Dominik Breitenbacher, and Yuval Elovici. N-baiot—network-based detection of iot botnet attacks using deep autoencoders. *IEEE Pervasive Computing*, 17(3):12–22, 2018.