

Stepping out of the MUD: Contextual threat information for IoT devices with manufacturer-provided behaviour profiles

Luca Morgese Zangrandi
luca.morgese@tno.nl
TNO
Den Haag, NL

Thijs van Ede
t.s.vanede@utwente.nl
University of Twente
Enschede, The Netherlands

Tim Booijs
tim.booijs@tno.nl
TNO
Den Haag, NL

Savio Sciancalepore
s.sciancalepore@tue.nl
Eindhoven University of Technology
Eindhoven, The Netherlands

Luca Allodi
l.allodi@tue.nl
Eindhoven University of Technology
Eindhoven, The Netherlands

Andrea Continella
a.continella@utwente.nl
University of Twente
Enschede, The Netherlands

ABSTRACT

Besides coming with unprecedented benefits, the Internet of Things (IoT) suffers deficits in security measures, leading to attacks increasing every year. In particular, network environments such as smart homes lack managed security capabilities to detect IoT-related attacks; IoT devices hosted therein are thus more easily infiltrated by threats. As such, context awareness on IoT infections is hard to achieve, preventing prompt response. In this work, we propose MUDSCOPE, an approach to monitor malicious network activities affecting IoT in real-world consumer environments. We leverage the recent Manufacturer Usage Description (MUD) specification, which defines networking whitelists for IoT devices in MUD profiles, to reflect consistent and necessarily-anomalous activities from smart things. Our approach characterizes this traffic and extracts signatures for given attacks. By analyzing attack signatures for multiple devices, we gather insights into emerging attack patterns. We evaluate our approach on both an existing dataset, and a new openly available dataset created for this research. We show that MUDSCOPE detects several attacks targeting IoT devices with an F1-score of 95.77% and correctly identifies signatures for specific attacks with an F1-score of 87.72%.

ACM Reference Format:

Luca Morgese Zangrandi, Thijs van Ede, Tim Booijs, Savio Sciancalepore, Luca Allodi, and Andrea Continella. 2022. Stepping out of the MUD: Contextual threat information for IoT devices with manufacturer-provided behaviour profiles. In *Proceedings of ACSAC '22: ACM Annual Computer Security Applications Conference (ACSAC '22)*. ACM, New York, NY, USA, 15 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

The Internet of Things (IoT) paradigm refers to a cyber-physical ecosystem of interconnected devices (things), which exchange and process data to enable intelligent decision-making [1]. IoT adoption is in continuous growth, as leading agencies forecast the number of Internet-connected devices to quadruple from 9 billion in 2020 to 38 billion in 2030 [2]. Unfortunately, increasing IoT adoption also leads to an increased attack surface, as recent works reveal. In 2020 only, the number of IoT-directed attacks increased by 800% compared to the previous year [3] — mostly driven by a proliferation of

post-Mirai botnets (e.g., Mozi, [4] OKIRU, [5] BrickerBot, [6] and Persirai [7]). A Netscout report [8] reveals that IoT devices are attacked “within the first five minutes of their deployment”. One and a half billion attacks on smart devices were recorded in the first half of 2021, more than twice than in the same period in 2020 [9]. At the same time, research shows that a considerable share of compromised IoT devices is located within consumer environments (Internet Service Providers (ISPs) and telecommunications internet domains [10, 11]), pointing to IoT environments that we refer to as loosely-protected: households, shops, restoration, open-access Wi-Fi areas, and alike. IoT attackers actively infiltrate in these environments [12] where compromised IoT can proliferate largely undetected by the defenders’ community.

To improve detection, research and industry stress the need to enhance continuous collection and sharing of actionable IoT security information, e.g., indicators of compromise, IoT vulnerabilities, attackers’ goals and trends [10, 13–16]. This information allows faster, more precise, and better-informed security interventions on both network administrators and device vendors’ sides [10, 17–19], making IoT integration and use safer.

In the state-of-art, there are three main approaches to collect and share security information: (IoT) honeypots, network-telescopes, and Threat Intelligence (TI) feeds. Each, though, has limitations with respect of visibility of real-world loosely-protected IoT environments. IoT honeypots mimic vulnerable devices to study how attackers interact with them [12, 20, 21]. However, they can be fingerprinted by attackers, and do not directly represent how anomalies spread through real-world consumer deployments [12, 22, 23]. Network telescopes (or, ‘darknets’) are portions of routable IP addresses that do not host any service; all traffic that they receive is thus anomalous by definition [24]. The insights network telescopes derive are biased towards Internet-wide activities [10, 25–27], i.e., no targeted IoT attacks can be detected [19]. Crowd-sourced threat intelligence feeds such as AlienVaultOTX, Censys, and Shodan [28–30], when IoT-specific, can be highly inconsistent in reporting timings and coverage of attacks [27, 31, 32]. Lastly, other approaches such as DShield [33], collect firewall logs from multiple deployments, and produce large-scale indicators on network compromise attempts. However, they are not IoT- nor deployment-specific, and require additional configuration at the deployments’ end, making them not suitable for an average consumer [34].

In summary, we identify a gap in the state-of-art, about accessing and leveraging threat intelligence on IoT network threats that target consumer environments. Motivated by this, this work presents MUDSCOPE, a tool implementing a novel approach to monitor malicious IoT traffic in consumer environments. MUDSCOPE characterizes the observed anomalies, and can provide information on how threats spread through geographically-distributed IoT deployments. Our approach can assert on a timely basis what specific devices and deployments are targeted by similar or different network anomalies, by generating traffic signatures, and comparing them for multiple devices. Similar anomalous traffic signatures represent similar malicious activities from a network threat.

Our approach leverages the recent IETF [35] Manufacturer Usage Description (MUD) specification [36]. MUD allows vendors to whitelist traffic for their devices, ensuring that any non-whitelisted traffic can safely be rejected. MUD is an effort to improve the security of IoT with a specification-based approach. Previous work has shown the effectiveness of MUD profiles to filter malicious IoT traffic [37–40]. In brief, all traffic that is not whitelisted in MUD profiles is thus necessarily anomalous.

We build on MUD profiles and their proved effectiveness, and use the specification in a novel way, to analyse MUD-rejected traffic (MRT), and generate previously-unavailable threat intelligence. For a specific device in a deployment, we monitor its MRT, and describe how it evolves over time: we cluster rejected traffic flows showing similar features, and we describe how these clusters evolve. The evolution of these clusters provides a signature for anomalous activities related to an IoT device. We compare signatures collected from multiple devices to gather insights on IoT network anomalies and how they spread through loosely-protected environments.

We validate MUDSCOPE with six IoT devices from different categories and brands. We target them with different scanning and Denial of Service (DoS) attacks, collect and generate signatures of MUD-rejected traffic for each device, and show that we can detect when similar or different network threats reach multiple devices. This information can be provided to manufacturers of targeted IoT products and deployments, and the research community, to prompt investigation and response.

In summary, we make the following contributions:

- We propose a new approach for monitoring IoT-related traffic that leverages the MUD specification to reveal malicious connections. This operates as an IoT-specialised network telescope requiring minimal configuration.
- We implement our approach in MUDSCOPE, a tool that generates attack signatures describing the behavior of clustered anomalous network flows over time.
- We evaluate MUDSCOPE on our new, open-source, dataset containing multiple IoT devices from different categories and brands. Our experiments show that MUDSCOPE can characterize malicious network traffic with an F1-score of 87.72%.

In the spirit of open science, we make MUDSCOPE open-source [41], as well as the data-sets generated and used for our evaluation [42].

2 THREAT MODEL

Cyber attackers are actively engaged in compromising exposed IoT devices [8, 9, 12]. Of the targeted devices, over a half are deployed

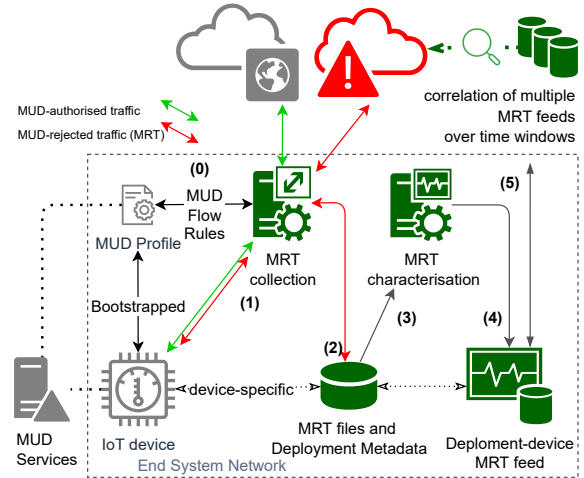


Figure 1: High-level schematic for the proposed MUD usage. (0) The local environment retrieves MUD profiles and flow rules; (1) we enforce MUD rules on device-related traffic; (2) we save the MUD-Rejected Traffic (MRT) with associated metadata; (3) we describe how the MRT evolves through time; (4) we save the description in an MRT feed; (5) we compare multiple MRT feeds to observe fluctuations for different IoT devices.

within ISPs and telecommunications sectors [10, 11]. This points to devices deployed in loosely protected environments such as households and home offices. These environments lack managed security capabilities such as intrusion detection systems and security operation centers. Thus, the devices that they host are more easily reached by network threats [43].

In this threat model, MUDSCOPE is designed to detect anomalous IoT network activities affecting devices in loosely-protected environments. We focus our research on consumer IoT devices (e.g., IP-cameras, motion sensors, smart appliances, smart plugs, etc.) that communicate over the Internet via UDP/TCP-IP stack. We address external network threats, such as botnets, that attempt to gain control of IoT devices via reconnaissance (active scanning) and initial access tactics, as outlined in the MITRE ATT&CK kill-chain [44].

The MUD specification is still not widely adopted by vendors, though active engagement from standardisation bodies, industry, and research is manifested [45]. In our research, we assume that MUD profiles are available, and readily deployable as OpenFlow rules [46] (forwarding rules for the allowed UDP/TCP flows associated to the MUD rules), as per indication of the MUD specification. To satisfy this assumption, we rely on MUDgee [37] to create MUD profiles. Finally, we assume that the integrity of deployed MUD profiles has been verified by a trusted component.

3 METHODOLOGY

A high-level scheme of our approach is presented in Figure 1. OpenFlow rules derived from a MUD profile are available in the local environment (step 0). In step 1, we collect device-specific traffic packets rejected by the MUD rules, namely, the MUD-Rejected Traffic (MRT). At step 2, we record MRT packets logs divided in time windows. In step 3, we describe anomalous traffic at each time

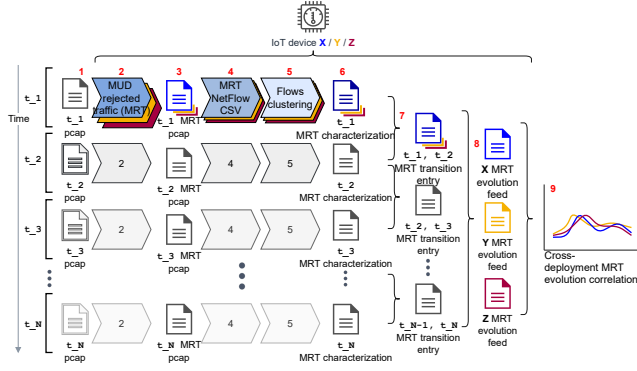


Figure 2: Pipeline for our MUD usage proposal. The sequence of steps of our pipeline is detailed for the time window t_1 . One colour codifies one device.

window: we collect traffic flows and cluster them together, obtaining a high-level representation of anomalies through that given time-window. At step 4, we create traffic signatures by describing how clusters of anomalous flows change through time windows. We refer to these description as MRT feeds. These feeds contain signatures of anomalous activities, for a specific device-deployment. Finally, at step 5, we compare signatures from multiple MRT feeds to detect emerging similar attacks patterns.

Note that all the steps in the pipeline are specific to one device in one deployment, up to the comparison of traffic signatures from multiple devices. We explain these steps referring to Figure 2, presenting the pipeline of our approach in greater detail.

3.1 MUD profiles enforcement

Our approach starts from MUD profiles available for devices within a loosely-protected environment. In particular, in accordance with the MUD specification, the whitelist rules in a MUD profile have been translated into actionable networking rules [36] – in our case, an OpenFlow table. OpenFlow is a widely adopted protocol to control packets forwarding in switches or routers in a software-defined (i.e., programmable) network [47] – a framework often coupled with IoT [48]. MUD rules are therefore enforced at local level via network flow rules.

Throughout a time window t_a , we listen to device-specific network packets, and test each packet against the MUD profiles. If no match is found, the packet is anomalous, and appended to a network log (a pcap) file. We then convert the rejected packets in t_a to netflow [49] flows format to constrain the size of the MRT logs. These operations correspond to steps 1, 2 and 3 in Figure 2.

3.2 Time-window anomalies characterisation

The MUD-Rejected Traffic is anomalous by manufacturer definition. In general, we expect that collected MRT flows belong to certain categories of anomalous traffic, such as host-discovery probes, targeted scans, credentials brute-forcing, DoS traffic, advertisement-related probes, internet noise [19]. Moreover, different threat actors, such as botnets, may attempt to compromise new hosts using signature routines [50, 51]. For instance, in its first infection phases, MIRAI targets a victim device via (1) scanning TCP ports 23 and 2323, (2)

brute-forcing Telnet credentials, (3) injecting code upon access to the device [52]. Our approach is primarily aimed at detecting when several devices are targeted by similar anomalies underlying similar attack routines.

We make two observations: (i) the same type of malicious activities, observed through *one* time window t_a , will likely yield similar MRT profile; (ii) the same attack routine will yield a similar sequence of anomalous activities observed through *multiple* time windows.

With the above observations, we adopt the following approach to group together flows belonging to similar malicious traffic. We first conducted a preliminary analysis on a dataset for IoT network intrusion detection (Kang et al. [53]) to select the flow features that most discern IoT network attacks. We concluded that the six flow features of (i) bytes-per-packet, (ii) TCP flags, (iii) input and (iv) output bytes, (v) destination port, and (vi) source-address category (private, public, reserved) are good indicators to monitor to discern different IoT attack types (we report this analysis in Section 6.2). We use these features to cluster together similar anomalous flows observed through a window t_a (step 5 in Figure 2). We choose HDBSCAN [54] as clustering algorithm, due to its time efficiency, and highly adaptive and generalising nature [54, 55], which suits the dynamic and unknown nature of network traffic [56].

We describe the ensemble of observed anomalous flows in t_a by means of the characteristics of the clusters space, i.e., the number of observed clusters and their positions. Specifically, we describe the clusters' positions with the list of their *meta-centroids*. A meta-centroid of a cluster is a point defined as follows: one dimension per average value of the flow feature of all points in the cluster; plus two *meta*-dimensions for average and standard deviation values of between-cluster points distances. A meta-centroid is thus $n+2$ -dimensional, with n = number of MRT flow features (six in our case). We add the last two dimensions to have a more robust notion of the identity of a cluster. For instance, in two different time windows, two clusters may appear in similarly-centered regions, but present flows with different features. The meta-dimensions allow us to capture this difference.

Finally, the MRT for a device observed through a window t_a is represented by the set of $(n+2)$ -dimensional clusters meta-centroid points $C_a: \{c_i^a: i=0, \dots, |C_a|\}$ – a *characterisation*, step 6 in Figure 2 – with c_i^a the meta-centroid of the cluster i , for the clustering performed in the time window a .

3.3 Anomaly signature

Thus far, we collected the MRT for a specific device, in a given time window. We then distinguished different clusters in this traffic. The idea behind this is that each attack is represented by a one or more clusters. We now want to describe, in turn, how the positions of the clusters change through successive time windows to track how an attack evolves. The track of the changes in the clusters space constitutes the *signature* of an anomalous event.

We consider the general case of two successive characterisation: $C_a: \{c_i^a, i=0, \dots, |C_a|\}$, over t_a , and $C_b: \{c_j^b, j=0, \dots, |C_b|\}$, over t_b . We assume that t_a precedes t_b , i.e., the end-time of t_a comes before the start-time of t_b : $t_{a_start} \leq t_{a_end} \leq t_{b_start} \leq t_{b_end}$.

We compute the distance matrix M between pairwise clusters, where $M[i, j] = \text{dist}(c_i^a, c_j^b)$, $i=0, \dots, |C_a|, j=0, \dots, |C_b|$. Note that row

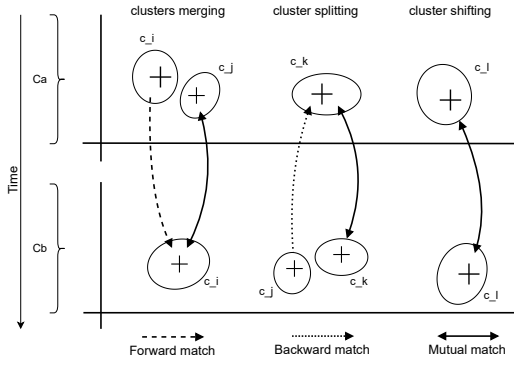


Figure 3: Conceptual illustration of changes in the anomaly clusters space through two time windows. The clusters space is simplified to two generic dimensions. We describe the dynamic behaviour of anomalies by means of mutual, forward, and backward match cases.

indices correspond to clusters in C_a , and column indices to clusters in C_b . We observe what clusters appear closest in space, in the two different time windows, and define three match cases that describe shifts, splits, and merges of clusters.

The cluster $c_i^a \in C_a$ **forward-matches** with the cluster $c_j^b \in C_b$ if c_i^a is closest to c_j^b among all clusters in the *future* characterisation C_b . Analytically, $m_{\text{fwd}}(c_i^a) = c_j^b \iff$

$$c_j^b = \operatorname{argmin}_{j=0, \dots, |C_b|} (\operatorname{dist}(c_i^a, c_j^b)).$$

Similarly, a cluster c_j^b **backward-matches** with c_i^a when the cluster c_j^b in the new time window is closest to the cluster c_i^a from the previous characterisation C_a . We say $m_{\text{bwd}}(c_j^b) = c_i^a \iff$

$$c_i^a = \operatorname{argmin}_{i=0, \dots, |C_a|} (\operatorname{dist}(c_i^a, c_j^b)).$$

A **mutual match** between c_i^a and c_j^b , for fixed i and j , occurs when

$$m_{\text{fwd}}(c_i^a) = c_j^b \wedge m_{\text{bwd}}(c_j^b) = c_i^a.$$

Forward matches are only defined from C_a to C_b , and backward matches are only defined from C_b to C_a .

Figure 3 illustrates the clusters space in two consecutive characterisations. Clusters c_i^a, c_j^a , merge in c_l^b . They trigger one forward-match, and one mutual-match. Cluster c_k^a splits in two clusters c_m^b, c_n^b , and triggers one mutual-match and one backward-match. Cluster c_i^a shifts to the position of c_l^b , triggering a mutual-match.

The two distance matrices in Table 1 show two examples of clusters evolution, generated with our implementation (see Section 4). In the first matrix, we see a relatively stable evolution: just one new cluster (c_3^b) appears at a very short distance from the existing ones. In the second matrix, the high agglomeration of forward matches highlights a *merge* of clusters to c_3^a, c_5^a , and c_6^a , to c_3^b . We illustrate how network anomalies relate to these descriptions of clusters evolution in Section 3.3.2.

$C_a \backslash C_b$	0	1	2	3	4	5
0	0.020858	2.419684	2.405222	2.403964	2.403957	2.403956
1	2.416821	0.001128	0.281784	0.277283	0.277189	0.277154
2	2.402623	0.281764	0.001103	0.052307	0.052008	0.051955
3	2.400945	0.277189	0.051946	0.003871	0.001101	0.002798
4	2.400944	0.277154	0.051893	0.006378	0.002798	0.001101

$C_a \backslash C_b$	0	1	2
0	3.019654	3.596038	3.825001
1	0.354194	2.373034	2.351585
2	2.495640	0.162561	0.352583
3	2.464904	0.412329	0.247349
4	2.399984	0.330156	0.099545
5	2.564678	0.501293	0.412108
6	2.476171	0.341645	0.176153

Table 1: Two scenarios exemplifying cases of mutual, forward, and backward matches between clusters over consecutive time windows, taken from our tool. In the case above, MUD-rejected traffic is relatively stable from C_a to C_b : only one backward match at a very close distance is recorded. In the case below, we observe a merge of clusters: three forward matches agglomerate over cluster c_3^b .

3.3.1 Signature definition. The three defined match cases allow us to derive features to describe the evolution of the MRT through consecutive characterisations C_a and C_b . We extract the following ten metrics: (1) balance of gained-to-lost clusters (clusters balance), to track the number of clusters observed over time; (2) average of all distances over the distance matrix of meta-centers (all dists avg), describing how spread clusters are; (3) mutual matches n , the number of mutual-match cases; similarly, (4) backward matches n and (5) forward matches n . We account for the respective percentage of match events of each type over all match events: (6) mutual matches percentage, (7) backward matches percentage, and (8) forward matches percentage, to record the relative share of each match type. Finally, (9) bwd matches agglomeration avg and (10) fwd matches agglomeration avg indicate over how many clusters on average forward and backward matches agglomerate, capturing the volume of splits and shifts of single clusters.

We synthesize these features with the assumption that they achieve a high-level description of a transition between two MRT cluster spaces. We evaluate their effectiveness in Section 6.1.

For consecutive pairwise characterisations of MRT, we compute and write these features to MRT *transition entries* (step 7 in Figure 2). Each entry describes how anomalous traffic changes through two time windows. The sequence of MRT transition entries constitute an MRT feed (step 8 in Figure 2), specified for the time range spanning from the start-time of the first window to the end-time of the last window.

According to this definition, the signature of an anomaly can be represented by a segment of an MRT feed. This is an $F = [N, M]$ matrix, where N is the number of MRT transition entries through which the anomaly spans, and M is the number of signature features considered in the signature. A column of $F[:, j]$ represents the values of the signature feature j through the time span where anomalous activity occurs. We denote the i^{th} signature in a feed X as F_i^X .

3.3.2 Expected behaviour. If a device is subject to anomalous but harmless *noise* [57], over time, we expect the values of the features (i.e., columns) in its MRT feed to show little variation. The number of clusters will remain stable, and the percentage of mutual matches will stay close to 100%.

A *new* network event would generate new clusters, increasing the clusters balance, and the backward matches. The backward match agglomeration will also increase, because previous clusters are expected to mutually match. If an anomaly *ceases*, then the number of clusters will reduce in the next characterisation, and forward matches will increase — corresponding to the previously existing clusters pointing ahead towards the remaining clusters.

3.4 Comparison of MRT feeds

A network threat using the same attacks to compromise different IoT devices will cause similar fluctuations in the MRT feeds of the respective devices. To detect this (step 9 in Figure 2), we perform the following three steps. First, we take the MRT feeds of the devices we are interested in monitoring, over arbitrary time periods. Second, we detect if and where anomalous traffic occurs in the feed, by observing if any MRT is captured in any time window in the first place. This is reflected in the signature feature reporting the clusters' balance through time windows. Doing so, we automatically gather the signature for any anomaly recorded in each feed. Third, we compare all anomaly signatures, pairwise. If the signatures length N is different, we compare the smaller signature with equally-sized sliding windows on the larger signature, and we record the values for the highest similarity found. In particular, for any two signatures F_n^X and F_m^Y , we compute the Pearson correlation coefficient r_j for each two column vectors of the same feature j , to observe if their values change in a similar way. The average value of r for all features is the *correlation* c of the two signatures, i.e., how similar is the trend of the anomalous activity captured between the two signatures. Analytically,

$$c(F_n^X, F_m^Y) = \frac{1}{M} \sum_{j \in \text{features}} (r_j = r(F_n^X[:,j], F_m^Y[:,j]))$$

, where M is the number of MRT feed signature features considered.

To capture the case where c returns a low value despite some signature features showing high correlation, we also record the maximum value of r over all features: $m(F_n^X, F_m^Y) = \max_{j \in \text{features}} r_j$. Additionally, to make our alerting method robust to comparing signatures from *different* attacks but with *similar* trends, we record the proportion p between ranges of the number of clusters generated by the two signatures, expressed as Ra_n^X and Ra_m^Y :

$$p(F_n^X, F_m^Y) = \frac{\min(Ra_n^X, Ra_m^Y)}{\max(Ra_n^X, Ra_m^Y)}.$$

We then use these values to compute the signature match value:

$$a(F_n^X, F_m^Y) = \text{avg}(c(F_n^X, F_m^Y), m(F_n^X, F_m^Y)) * \xi,$$

where ξ is a slack variable that constrains the value of the alert based on $p(F_n^X, F_m^Y)$, and is defined as follows:

$$\xi = p(F_n^X, F_m^Y) \text{ if } p(F_n^X, F_m^Y) \leq 0.5, \text{ else } 1.$$

A match alert is prompted when $a(F_n^X, F_m^Y)$ is greater than a given threshold, chosen to be 0.5, to capture the case where only

one signature features correlates perfectly. Further approaches can be investigated in follow-up work.

4 IMPLEMENTATION

For our proof-of-concept, we generate MUD profiles and related OpenFlow rules with open-source software from previous research, i.e., MUDgee [37]. We collect local network traffic using Wireshark [58]. We divide such collected traffic in time-window pcaps, according to our methodology, using tcpdump [59]. These pcaps are then fed to MUDSCOPE, our tool. MUDSCOPE is implemented in Python 3.8. We make use of Python's Scapy [60], a network packet manipulation library, to match device-specific packets against MUD flow rules, and generate the MRT pcap files. Packets files are aggregated into network flows in Comma-Separated Value (CSV) files, using nfdump [61].

From a MRT flow file, we select the chosen attack-discriminating flow features (see Sections 3.2 and 6.2), which we pre-process with the Sklearn library [62] to quantify and scale to apply clustering. We use the HDBSCAN python implementation [54] to produce the clustering characterisation for each MRT flow file, as explained in Section 3.2.

For consecutive characterisations, we compute distance matrices among clusters' meta-centroids with Numpy's [63] euclidean distance function. From these matrices, we extract the MRT evolution signatures indicators presented in Section 3.3.1, mapping the sequential change of characterisation in a MRT transition entry. We build a MRT feed CSV file via appending consecutive such entries.

Finally, we implement a MRT feed monitor module that ingests an arbitrary number of feeds, extracts anomaly signatures, and compares the similarities between signature metrics using Numpy's Pearson correlation index `corrcoef`. We recall that MUDSCOPE is provided as open-source [41].

5 DATASETS

We evaluate our approach on both the existing IoT Network Intrusion Dataset by Kang et al. [53] and on our openly available MUDSCOPE IoT DATASET [42].

IoT Network Intrusion Dataset. Kang et al. [53] provide a dataset containing 42 raw network traces spanning 9 attacks and benign traces (2.99M packets) of two different smart home devices (a *Nugu* smart speaker and gan *EZVIZ* wi-fi camera) captured over a (non-continuous) period of 112 days. The attacks include scanning (host, port and OS) using Nmap; flooding (SYN, UDP, ACK and HTTP) using custom scripts and MIRAI botnet attacks; and Host discovery and Telnet brute-force attacks using the MIRAI botnet. The original dataset also contains MITM ARP spoofing attacks, but these attacks are considered out of scope for this research. The IoT Network Intrusion Dataset is used for parameter optimization of our method (Sections 6.2 and 6.3). Appendix B gives an overview of the dataset.

MUDSCOPE IoT DATASET. We also contribute our own dataset (see Table 2). This dataset was generated by attacking heterogeneous IoT devices deployed at a physical location L1 (anonymized for peer-review). The deployed devices are: Eufy security homeBase 2 doorbell, Honeywell thermostat, two different Foscam IP cameras (models C1780P and RM2), and two Hombli smart plugs. We deploy

Attack	Duration	# Flows	# Packets	# Devices
None (Benign)	1,462.93 s	627	76,055	5
Telnet/SSH port scan	1,663.82 s	424	4,815	5
OS scan	2,854.25 s	10,145	24,371	5
Vulnerability scan	3,815.83 s	3,091	14,801	5
TCP SYN flood DoS	2,131.20 s	135,472	395,771	2

Table 2: Details of MUDSCOPE IoT DATASET. For each attack scenario, we specify the capture duration, total number of flows, packets, and number of involved devices.

Device type	No. devices	Scans			DoS SYN flood
		Telnet/SSH	OS	Vulnerabilities	
Eufy HomeBase 2 doorbell	1	×	×	×	
Honeywell Round T57RF2025 thermostat	1	×	×	×	
Hombli smart plug HBPP-0201	2	×	×	×	
Foscam C1780P IP camera	1	×	×	×	×
Foscam RM2 IP camera	1				×

Table 3: Description of the devices utilised for the data collection, and the attacks performed on them.

an attacking computer in a different subnet at a second geographical location L2. The dataset consists of raw network traces for the six IoT devices containing both benign data and scanning attacks (Telnet SSH port scan, OS discovery and Vulnerability scan) as well as targeted (hence only two devices) DoS TCP SYN-flood attacks. Table 3 summarizes the devices and attack scenarios recorded in the dataset. Section 6.1.1 describes the procedure we used to generate the dataset. By capturing these various attacks over multiple devices in various locations and from different vendors, we attempt to show that our approach can correlate similar attacks in differing scenarios (Section 6.1).

6 EVALUATION

The main objective of our approach is to detect attacks on IoT devices and create respective signatures that can be used to identify emerging common attacks. Tables 4, 5, and 6 give an overview of our main results for detection and signature comparison (details in Section 6.1) over our MUDSCOPE IoT DATASET (Section 5).

To obtain these results, we first used the Kang dataset to systematically analyze which features are most relevant for generating signatures (Section 6.2), and to perform an intermediate analysis on the clusters generated by our approach (Section 6.3).

In this Section, we first report the main validation for MUDSCOPE in 6.1. We then report the analyses performed to develop our tool, in Sections 6.2 and 6.3.

6.1 Signature matching

We aim to ascertain whether and when multiple IoT devices have been targeted by the same network threat. To this end, MUDSCOPE correlates signature fluctuations of the MRT. To evaluate the extent in which these signatures are correctly correlated, we perform several experiments on six IoT devices. The network traces for the experiment we run constitute our MUDSCOPE IoT DATASET, as we describe in Section 5.

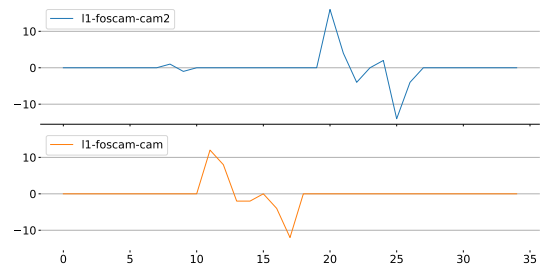


Figure 4: Anomaly signature fluctuations for the clusters balance feature, for the two Foscam cameras attacked with DoS traffic. The Y axes reports the ranges of the clusters balance values. The X axis refers to the MRT transition entries.

6.1.1 Experimental setup. We performed four experiments. In each of them, we start a packet capture at L1, then target a subset of the devices asynchronously with a given attack from L2, and stop the capture. The attacks that we launched are the following: (1) scan of Telnet and SSH ports, (2) nmap OS discovery scan, (3) a slow¹ nmap vulnerabilities testing scan, and (4) TCP SYN flood DoS. For the three scanning experiments, we left the RM2 Foscam camera unharmed to have a control recording. In the targeted DoS experiment, we attacked both Foscam cameras, and left the other devices unharmed.

We obtain one general capture per experiment, and process each with MUDSCOPE, for each device’s MUD profile. We so obtain four MRT feeds per device – one per experiment. We expect that each MRT feed contains the signature of the attack performed in the relative experiment, as generated by the MRT for the related device. Next, for each experiment, we extract the signatures in each related MRT feed for each device, and verify whether the malicious traffic was captured. Table 4 shows to what extent our detection approach finds signatures of attacks (anomalies) in the MUDSCOPE IoT DATASET. As shown, we correctly identified all anomalies underlying attacks (100.00% recall) and incorrectly identified 9 anomalies (91.89% precision).

Next, we want to validate that signatures from the same attacks match with each other, but not with signatures from different attacks. Therefore, we run two sets of experiments. First, we check whether the signatures from the same attack correlate among all attacked devices. We report the overall results of this validation in Table 5. There, we see that although there is a complete overlap of all signatures for simple attacks, such as the Telnet/SSH port scan, more device-interactive attacks (e.g., OS scan) report some False Negatives (we explain below in Section 6.1.2).

Second, to validate that signatures from different attacks do not match, we perform six tests as follows: for each of five devices (all except the Foscam RM2), we compare all of its MRT feeds containing an attack. In the sixth test, we compare the MRT feeds from four different devices (one Hombli plug, the Honeywell thermostat, the Eufy HomeBase, and the Foscam RM2), each with a different attack. We report the overall results of this second validation in Table 6. We publish the complete set of outputs of our results at [64].

¹We used the T_2 evasion timing template for the nmap operations. Refer to <https://nmap.org/book/performance-timing-templates.html>.

Attack	MRT entries	Anomalous entries	TP	TN	FP	FN	Accuracy	Precision	Recall	F1-score
Telnet/SSH port scan	130	15	15	115	0	0	100.00%	100.00%	100.00%	100.00%
OS scan	217	26	26	191	0	0	100.00%	100.00%	100.00%	100.00%
Vulnerability scan	310	48	45	259	3	0	98.06%	93.75%	100.00%	96.77%
TCP SYN flood DoS	170	22	16	142	6	0	96.47%	72.73%	100.00%	84.21%
Total	827	111	102	707	9	0	98.90%	91.89%	100.00%	95.77%

Table 4: Performance of anomaly detector over Novel IoT Dataset. We report the performance in terms of number of True/False Positive/Negative detected MRT entries as well as the Accuracy, Precision, Recall and F1-score of our anomaly detector.

Attack	Total signatures	Expected matches	TP	TN*	FP	FN	Accuracy	Precision	Recall	F1-score
Telnet/SSH port scan	5	10	10	N/A*	0	0	100.00%	100.00%	100.00%	100.00%
OS scan	5	10	6	N/A*	0	4	60.00%	100.00%	60.00%	75.00%
Vulnerability scan	5	10	8	N/A*	0	2	80.00%	100.00%	80.00%	88.89%
TCP SYN flood DoS	2	1	1	N/A*	1	0	50.00%	50.00%	100.00%	66.67%
Total	17	31	25	N/A*	1	6	78.13%	96.15%	80.65%	87.72%

*This experiment compared found anomalies only of the same attack type, hence the True Negatives are 0. Table 6 shows the results when comparing signatures of different attacks.

Table 5: Evaluation matching attack signatures. We show for each attack, whether our signature matching algorithm found a match between signatures of the same attack.

Test	Device(s)	Compared MRT feeds	Incorrect matches		
			Expected	Worst	Found
1	Eufy home-kit doorbell	Scans (Telnet/SSH, OS and Vulnerabilities)	0	3	0
2	Honeywell thermostat		0	3	0
3	Hombli plug 1		0	3	0
4	Hombli plug 2		0	3	0
5	Foscam camera C1780P	All	0	6	1
6	Eufy, Honeywell, Hombli, Foscam	All	0	6	0
Overall			0	18	1
Matches correctly discarded			94.44%		

Table 6: Summary of the tests executed for the signature-matching for different attacks. We expect that at each test, no matches are found, because we are comparing signatures from different attacks. In the worst case, we find that all signatures match with one another. For these tests, our method proved effective in not raising match alerts for different attacks.

6.1.2 Results for same-attacks signatures matching. As can be seen in Table 6, we detected all Telnet and SSH ports scans directed to each device, and found correlations between all signatures. In more detail, the attack consists only of four flows, directed to ports 22, 23, 2222 and 2323. For this reason, the produced fluctuations in the MRT feeds are minimal, and we are able to detect signature matches only by means of the `clusters` balance metric. The metric shows perfect correlation for this attack, while values for the other nine metrics in the extracted signatures are flat, and thus record a correlation value of 0. The RM2 Foscam camera, which was not attacked, does not record anomalies, as expected.

For the OS scan attack, we detect all anomalous flows, and record false-positive anomalous flows, simply because of benign packets that were accounted in the generated MUD profiles. We find all signature matches between all devices except the Foscam camera C1780P, thus returning six matches instead of ten. This occurs because, differently from the other four devices, the Foscam camera

engaged minimally with the OS scan, and thus its MUD-rejected traffic produced a different signature. This provides the insight that anomaly signatures for the same attack can vary depending on how the devices react, and can thus be device-specific.

In the slow vulnerability scan experiment, all anomalies are detected correctly, plus false positive anomalies for the same reason as above. We verified that the false-positive anomalies are simply generated from benign packets that the MUDgee tool failed to include in the generated MUD profiles. We detect eight out of ten matches, leaving out the match among Eufy and Honeywell, and Honeywell and one of the two Hombli plugs. Interestingly, all matching signatures correlate poorly (with an overall average of 0.33), but we are still able to detect the matches with our method, because of high maximum correlation values (0.9 on average), and comparable magnitude of generated clusters.

We detected both DoS attacks in the fourth experiment. We verified that the recorded false-positives are due to benign packets that were not included in the MUD profile, and barely generate a match on their signatures. Most importantly, the signatures of the two DoS attacks are very strongly correlated over all signature features (average of 0.86). In Figure 4, we report the plot for the fluctuations of the `clusters` balance signature feature, detailed for the the two Foscam cameras. It is possible to visually appreciate the similarity between the two plots. For this experiment, we report the complete plot for the `clusters` balance in in Appendix E, together with the reported anomalies correlation values over all features.

6.1.3 Results for different-attacks signatures matching. As shown in Table 6, as for the tests performed, our method for signature matching effectively discerns the different attacks, and does not raise match alerts in the very most of the cases. The only one false-positive match, reported among Foscam camera’s signatures, regards the Telnet SSH scan, and a false-positive anomaly recorded

Attack label	No. rejected flows	rejected % of total flows
dos-synflood	44,547	63.19%
unknown	10,413	14.76%
mirai-ackflood	7,659	10.85%
scan	5,815	8.24%
mirai-httpflood	1,737	2.47%
mirai-udpflood	223	< 0.01%
mirai-brutefrc-atk	94	< 0.01%
mirai-brutefrc-vict	59	< 0.01%
Total	70,547	-

Table 7: Layout of custom MRT flows dataset derived from Ezviz attack scenarios in Kang et al. [53].

Feature	Description	AMI score
bpp	bytes per packet	0.630
flgs_int	int value of flags bits array	0.585
da	destination address	0.518
oby	output bytes	0.468
ibyt	input bytes	0.456
dp	destination port	0.414
opkt	output packets	0.456
sa	source address	0.387
ipkt	input packets	0.307

Table 8: Selected features and their AMI scores with respect to the attack label.

during the vulnerability scan experiment, as described above. The average correlation value for the false-positive matching signature features scores relatively low, i.e., 0.13. This suggests again that, at least in terms of correlation, the two signatures are different. Though, the match is recorded because of similar magnitude of produced clusters, and a perfect correlation on the clusters balance feature, for those minimal fluctuations.

6.2 Flow Feature Selection

We use The *IoT Network Intrusion Dataset* by Kang et al. [53] to understand what flow features help the most in discerning different IoT network attack types. We first generate a MUD profile for the *Ezviz* camera using the benign traffic, and filter all anomalous traffic from the attack scenarios. Besides filtering most malicious packets (99.7%, with 91.3% in average for each attack scenario), noise traffic (i.e., local network traffic, and packets from Amazon, Google, Microsoft and alike) is also filtered, to which we assign the label *unknown*.

Appendix C displays the outcome of this step. We merge the MRT from all scenarios, and convert it to flows, labelled according to the attack they implement. The resulting labelled dataset is presented in Table 7. Each flow is described via the set of NetFlow features reported in Appendix A. We compute the Adjusted Mutual Information (AMI) score [65] of each flow feature with respect to the attack labels. AMI is a robust metric in the presence of unbalanced classes, as is the case in the dataset. Considering two variables, AMI measures how the entropy of one variable X is reduced once the value of other variable Y is known. A high AMI score means that X is useful in determining the value of Y . We report the computation of AMI in Appendix F.

Table 8 presents a list of top-scoring features and their description. It is worth noting that we resolved IP addresses to a category label among *private*, *public*, *reserved*, *broadcast*, not to bias the AMI score with specific IPs. From these features, we discard: *opkt* and *ipkt*, because they are necessarily correlated with the higher-scoring *bpp*, *ibyt*, and *oby*. Besides, we choose *sa* instead of *da* because we are more interested in discerning anomalous traffic based on different sources, rather than targets.

6.3 Clusters analysis

Using the features from Section 6.2, we cluster flows together. Here, the intuition is that flows related to an attack present similar characteristics and can therefore be grouped into clusters. We can evaluate to what extent our mechanism creates clusters containing the same attack (class) using the homogeneity [66] score h . An homogeneity score of 1 indicates that, for each cluster, all samples belong to the same class; vice versa, $h=0$ when all clusters only contain samples of different classes. Besides computing the homogeneity score, we want to minimize the number of produced clusters in relation to the present anomalies. The reason is that one flow per cluster would discern all anomalies and therefore give a perfect homogeneity score of 1, but will not be meaningful as a cluster. To illustrate, an incoming distributed DoS produces a high number of distinct flows, but it should ideally yield a single, or at least a small number of clusters.

We use the MRT flows dataset for the *Ezviz* device (Table 7) to perform a grid-search on HDBSCAN’s standard parameters *min_cluster_size* and *min_samples* [54] to find a local maximum for homogeneity, and minimum for the amount of yielded clusters. We achieve the best results ($h=0.866$ and 11 clusters over the eight attack labels) when *min_cluster_size* = 1.2% the size of the dataset, and *min_samples* = 0.2% the size of *min_cluster_size*. On our dataset, this configuration also produces just 0.04% noise points. Table 9 lists the clustering results with respect to the produced clusters.

As we show in Table 10, the very most of the flows (93.18% in average) from all attacks except *mirai bruteforce* and *UDP flooding* scenarios are part of clusters where their share is most represented (as per Table 9). This happens because, when running the clustering on the whole dataset, *min_cluster_size*’s value was selected greater than the size of bruteforce and udp-flooding classes, which are therefore undetected. We thus also run the clustering on a subset of the dataset containing only bruteforce and udp-flooding scenarios, and we observe that the clustering produces usable results in this case well. We report these results in Appendix D.

Overall, these results show that our tool effectively groups anomalous flows of the same type in distinct clusters with high homogeneity. The sensitivity to different anomalous events yields a descriptive set of clusters to characterise the MRT at each time window.

7 DISCUSSION

On the signature-matching results. Overall, for the experiments that we performed, our method was able to discern 96.15% of the cases when two anomalies were from the same *type* of attack, and 94.44% of the cases when anomalies were from different types of attacks. In particular, it is worth noting that our experimental procedure included network attacks to devices ranging from four packets (the Telnet SSH ports scan), to tens of thousands of packets (the DoS

Cluster ID	Cluster size in % MRT flows	Most represented attack in cluster	Share of most represented attack
0	2.27%	dos-synflood	94.90%
1	11.00%	unkown	100.00%
2	1.22%	unknown	100.00%
3	1.66%	dos-synflood	99.91%
4	2.43%	dos-synflood	100.00%
5	56.16%	dos-synflood	99.99%
6	10.87%	mirai-ackflood	99.85%
7	1.26%	mirai-httpflood	97.42%
8	2.43%	scan	78.08%
9	5.85%	scan	96.87%
-1	4.79%	-	-

Table 9: Clustering results over produced clusters. Noise points map to cluster -1.

Attack label	Total	Represented in a cluster	
		No. flows	% of flows
dos-synflood	44,547	44,038	98.85%
unknown	10,413	8,632	82.90%
mirai-ackflood	7,659	7,659	100.00%
scan	5,815	5,342	91.86%
mirai-httpflood	1,737	871	50.14%
mirai-udplood	223	0	0.00%
mirai-brutefrc-atk	94	0	0.00%
mirai-brutefrc-vict	59	0	0.00%

Table 10: Percentage of flows that are part of a cluster where their class is the most represented.

attacks), as well as a more time-spread scan (the slow vulnerability scan). Therefore, our approach proved to be applicable to different types of network anomalies.

Another interesting aspect is the amount of data generated by our approach. After all, if a large organization or manufacturer monitors many devices, storage may become a problem. As we report in Table 5, the amount of data we generate by processing a capture of 30 minutes for 6 devices (5 attacked, one unharmed), consists of 310 entries, each recording a limited amount of data (i.e., practically the anomaly signature features, plus start and end timestamps), for time windows of 30 seconds each. We store 8 bytes per column value, with 20 columns per transition.

Therefore, considering 2 MRT feed entries produced per minute, per device, an MRT feed for 24 hours divided in time-windows of 30 seconds would yield a file of 27.65 MB size, which is reasonable to store in the short term. Now, MRT feeds can be removed when an attack has been investigated. Moreover, MRT feeds that do not capture anomalies do not need to be stored at all. On top of this, it would be practical to only store anomaly signatures, yielding even smaller files. Finally, incoming anomalies could be matched online against stored signatures databases. While these file sizes are manageable for smaller settings, we believe it is worth investigating the scalability of this approach to larger-scale scenarios.

On the methodology. In the methodology that this work presents, we cluster rejected flows in the attempt of capturing a high-level and concise description of an anomaly. Though we achieved this

in the clustering evaluation, in our experimental evaluation we observed that different anomalies produce varying numbers of clusters — i.e., the (more volumetric) attacks we launched were not represented by a single cluster, but instead by multiple clusters.

Though this performance is somewhat unexpected, it does not affect the effectiveness of the approach, because we define a signature precisely through the evolution of the clusters' space (recall Section 3.3.1). Besides, we showed that indeed we obtain signatures corresponding to different attacks.

Because this approach captures the notion of similar and different attacks, it gathers additional and novel situational awareness regarding network threats affecting consumer IoT devices. Furthermore, we note that the proposed clustering-based signature generation methodology outputs anonymized data feeds, which can be consumed by interested parties without privacy concerns.

On MUD adoption. As mentioned in Section 2, MUD profiles are not yet widely adopted by IoT manufacturers, although research and standardising bodies start to adopt MUD [67]. At the very least, MUD contributes to specification-based security for IoT devices.

In turn, one relevant consideration with respect to this work is the additional benefit that IoT manufacturers could gain from adopting MUD, and integrating an approach alike to ours. By adopting our solution, IoT vendors could monitor the intensity and coordination of malicious activities directed to their products, in a way that preserves the privacy of their customers. Thus, they would be able to detect when a particular model receives anomalous traffic from the same attack type at various customers. Besides suggesting that the device might present unknown vulnerabilities, such threat intelligence would provide insights on attackers' interests and trends, and promptly reveal the emergence of large-scale malicious events, such as new or specific botnets.

We finally note that our approach merely needs the NetFlow rules derived from MUD profiles in order to be operational. In general, our results promote the use of MUD profiles as an instrument to ascertain threat situational awareness for IoT devices. Thus, we believe that these results motivate further research in threat-intelligence collection methodology that we propose, and the exploration of others alike.

7.1 Limitations and future work

As our approach relies on MUD profiles for capturing network anomalies, it inherits their limitations. MUDSCOPE is not able to capture network attacks that can evade MUD rules, such as vendor compromise, man-in-the-middle, or spoofing attacks.

With respect to our signature-matching methodology, we observed some limitations regarding the attacks that we can effectively isolate and match. Anomalies producing very small fluctuations in the MRT feed values, such as a Telnet ports scan, or one probe packet from internet-wide scans, may all produce minimal and highly similar signatures. This would make it hard to discern these anomalies. Future work should investigate more advanced methods to differentiate these anomalies. One method could regard adding a signature feature recording the distance from the origin of the clusters space, to observe where small clusters appear, as related to the features of their flows.

Additionally, it may be hard to detect matches for attacks that are spread over a long period of time, as they could generate intermittent small signatures. Future work could develop a higher-level alerting module, performing comparisons on sequences of signatures.

With the hypothesis of a MUDSCOPE-aware attacker, we also note that our signature-matching methodology could be circumvented, e.g., by injecting noise traffic to a network attack, or mimicking signatures of non-harmful or general anomalies, such as internet-wide scans from known entities. This would still lead to a detection of the attack, but its signature will be unknown.

We finally account that we could not validate MUDSCOPE with authentic MUD profiles, and we cannot therefore grant that the profiles we generated are a perfect representation of authentic profiles.

8 RELATED WORK

IoT threat landscape monitoring. We identify in IoT honeypots and network telescopes the two state-of-art approaches that mostly align with MUDSCOPE’s objective of monitoring the IoT network threats landscape. We review some relevant works in the following.

One widely appreciated contribution to IoT honeypots is IoTpot, from Pa et al. [68] in 2015. Thanks to an adaptable backend, IoTpot studies Telnet-based IoT attacks directed to 8 different CPU architectures. The authors contribute with an analysis on the ‘scope and variety’ of IoT Telnet attacks. A 2018 work from Kamoen [20] builds upon the IoTpot custom-backend approach to further propose a persistent-state IoT honeypot (‘Honeytrack’). Therein, they maintain the status of the attack progress from a threat agent, and they re-presented it to the agent upon their successive interactions with the honeypot. Kamoen’s Honeytrack focuses specifically on studying the adversary behaviour, and how it dynamically evolves when a target is compromised and weaponised. A 2020 work from Tabari and Ou [12] similarly addresses the challenge of “largely unknown nature of attackers’ activities towards IoT”, by proposing a honeypot whose interaction is incrementally designed and integrated. They do so to progressively understand attackers’ specific interests, and thus interface simulated devices with higher-chance of compromise. They deploy their honeypot in 12 worldwide-spanning locations. Another 2020 contribution, by Wu et al., [21] proposes a controller architecture (‘ThingGate’) to broker configuration and communication data for bare-metal IoT honeypots. Their work is motivated by the need of studying IoT attacks in greater detail through physical honeypots.

IoTpot and Honeytrack’s authors make use of darknets to gain preliminary results suggesting what to account for in honeypot architectures. Indeed, darknet-based approaches are able to produce at-scale insights on IoT threats. A 2018 work from Shaik et al. [26] achieves internet-scale monitoring of compromised IoT devices, by correlating network telescope captures with threat intelligence feeds. Pour et al. [18] integrate the same setup with geolocation databases and ISPs’ feedbacks, and infer at-scale IoT-probing campaign characterized both by affected industry sector and vendors. The same research group expands the approach to achieve at-scale and locality-specific IoT-botnets evolution [25] and consumer-IoT compromises [27]. Furthermore, they implement ‘ex-IoT’, an IoT threat intelligence feed that streams findings from such network-telescope internet-scale IoT monitoring capability [69]. Notably,

a 2019 work from Griffioen and Doerr [51] studies Mirai-like botnets evolution and behaviour by leveraging on 7,500 Honeytrack [20] deployments, Delft’s University network telescope, and flow probing of infected devices.

Differently from the above approaches, our method proposes using real IoT deployments as a vantage point to collect malicious traffic. By design, this offers greater scalability options, and an upfront position to intercept malicious phenomena proactively. Finally, our work allows differentiating malicious traffic according to deployment characteristics, achieving a highly specific view on what are the targets of emerging anomalies.

MUD specification. The state-of-art on MUD profiles mostly concerns studies on its effectiveness as threat prevention tool, and proposals to extend their functionality.

One first contribution to MUD research is a 2018 work by Hamza et al. [37]. The authors create a tool to generate MUD profiles from network capture files of benign traffic (*MUDgee*). Studying the MUD profiles for 28 devices, they show how the specification can be integrated with Supervisory Control and Data Acquisition (SCADA) policies. A 2018 work from Schutijser [70] also proposes a MUD-profile generation tool, and shows how the specification is effective in blocking DoS attempts. Hamza et al. build upon MUDgee and study the intrusion-detection effectiveness of the specification [38]. They show that MUD is able to block internet-side threats, for devices with limited functionality, and specifically against volumetric attacks [39]. The researchers extend their MUD-based intrusion detection method in a work from 2019 [40], where they infer different anomalous status cases of devices, by matching their dynamically observed behaviour against their MUD-expected behaviour.

Other works are focused on MUD integration and enhancing profiles’ functionalities. Matheu et al. propose a way to extend the profiles to integrate security-testing results [71]. In related work, the authors design an SDN-backed authentication messages exchange to bootstrap MUD profiles in industrial environments [72]. Sajjad et al. [73] extend MUD profiles with firmware integrity information, fetched from vulnerabilities repositories, and distributed through a blockchain framework. Furthermore, they design a MUD bootstrapping routine that also accounts for gateway authentication, to prevent attackers to bypass MUD through router vulnerabilities. Ferardo et al. [74] propose a federated-learning framework [75] where only MUD-compliant devices are allowed to publish training data. Finally, Afek et al. [76] implement an ISP-level service to enforce MUD on behalf of SOHOs, unburdening deployments from the task.

Notably, the work of Afek concludes by stating that such ISP-level MUD services are at a perfect vantage point for detecting ‘global phenomena’ of malicious events affecting SOHO IoT (§VI, Afek et al. [76]). In fact, to the best of our knowledge, no other MUD work has yet explored this aspect. We acknowledge the state-of-art results on MUD as an intrusion detection tool, and as a base to instrument specification-based IoT security. We move from these findings to propose a novel use-case of MUD, focused on analysing the traffic that is rejected by the profile. Doing so, we gather insights on network threats targeting consumer IoT devices.

9 CONCLUSION

In this work, we presented an approach to gather insights into network threats targeting IoT devices. We showed that we can produce attack signatures that identify attacks on IoT devices from various manufacturers. We argued that this technique can be leveraged by both IoT manufacturers and the defenders' community to aid in the prompt detection and analysis of emerging IoT threats. Our approach is based on the MUD specification and leverages the advantages of specification-based IoT security.

We implemented our approach and released it as an open-source tool, MUDSCOPE, and we validated its performance on both an existing dataset and our own, openly available dataset. We showed that MUDSCOPE detects attacks with an F1-score of 95.77% and correctly identify signatures to a specific attack with an F1-score of 87.72%.

REFERENCES

- [1] ENISA. Internet of things (IoT) – ENISA. [Online]. Available: <https://www.enisa.europa.eu/topics/iot-and-smart-infrastructures/iot>
- [2] A. Holst. IoT connected devices worldwide 2019-2030. [Online]. Available: <https://www.statista.com/statistics/1183457/iot-connected-devices-worldwide/>
- [3] D. McMillen. Internet of threats: IoT botnets drive surge in network attacks. [Online]. Available: <https://securityintelligence.com/posts/internet-of-threats-iot-botnets-network-attacks/>
- [4] D. McMillen, W. Gao, and C. DeBeck. A new botnet attack just mozied into town. [Online]. Available: <https://securityintelligence.com/posts/botnet-attack-moziozied-into-town/>
- [5] Check Point Research. (2017-12-21) Huawei home routers in botnet recruitment. [Online]. Available: <https://research.checkpoint.com/2017/good-zero-day-skiddie/>
- [6] Trend-Micro. BrickerBot malware emerges, permanently bricks IoT devices - security news. [Online]. Available: <https://www.trendmicro.com/vinfo/us/security/news/internet-of-things/brickerbot-malware-permanently-bricks-iot-devices>
- [7] Tim Yeh, Dove Chiu, and Kenney Lu. Persirai: New IoT botnet targets IP cameras. [Online]. Available: https://www.trendmicro.com/en_us/research/17/e/persirai-new-internet-things-iot-botnet-targets-ip-cameras.html
- [8] NETSCOUT. (2021) Cyber security & threat intelligence report | NETSCOUT. [Online]. Available: <https://www.netscout.com/threatreport/>
- [9] T. Seals. (2021-09-06) IoT attacks skyrocket, doubling in 6 months. [Online]. Available: <https://threatpost.com/iot-attacks-doubling/169224/>
- [10] N. Neshenko, E. Bou-Harb, J. Crichigno, G. Kaddoum, and N. Ghani. "Demystifying IoT security: An exhaustive survey on IoT vulnerabilities and a first empirical look on internet-scale IoT exploitations," *IEEE Commun. Surv. Tutorials*, vol. 21, no. 3, pp. 2702–2733, 2019. [Online]. Available: <https://doi.org/10.1109/COMST.2019.2910750>
- [11] S. Torabi, E. Bou-Harb, C. Assi, M. Galluscio, A. Boukhtouta, and M. Debbabi. "Inferring, characterizing, and investigating internet-scale malicious IoT device activities: A network telescope perspective," in *48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks, DSN 2018, Luxembourg City, Luxembourg, June 25-28, 2018*. IEEE Computer Society, 2018, pp. 562–573. [Online]. Available: <https://doi.org/10.1109/DSN.2018.00064>
- [12] A. Z. Tabari and X. Ou. "A first step towards understanding real-world attacks on IoT devices," *CoRR*, vol. abs/2003.01218, 2020. [Online]. Available: <https://arxiv.org/abs/2003.01218>
- [13] M. van Staaldnuin and Y. Joshi. "The IoT security landscape: adoption and harmonisation of security solutions for the internet of things," TNO, Tech. Rep., 2019. [Online]. Available: <https://repository.tno.nl/islandora/object/uuid%3A989e7450-206f-4f7c-93aa-5587e4674781>
- [14] J. Saleem, M. Hammoudeh, U. Raza, B. Adebisi, and R. Ande. "IoT standardisation: challenges, perspectives and solution," in *Proceedings of the 2nd International Conference on Future Networks and Distributed Systems*, ser. ICFNDS '18. Association for Computing Machinery, 2018-06-26, pp. 1–9. [Online]. Available: <https://doi.org/10.1145/3231053.3231103>
- [15] O. Garcia-Morchon, S. Kumar, and M. Sethi. (2019-04) RFC 8576 - internet of things (IoT) security: State of the art and challenges. [Online]. Available: <https://tools.ietf.org/html/rfc8576>
- [16] A. Costin and J. Zaddach. "IoT malware : Comprehensive survey , analysis framework and case studies," in *BlackHat USA*, 2018. [Online]. Available: <https://i.blackhat.com/us-18/Thu-August-9/us-18-Costin-Zaddach-IoT-Malware-Comprehensive-Survey-Analysis-Framework-and-Case-Studies-wp.pdf>
- [17] M. Husák, N. Neshenko, M. S. Pour, E. Bou-Harb, and P. Čeleda. "Assessing internet-wide cyber situational awareness of critical sectors," in *Proceedings of the 13th International Conference on Availability, Reliability and Security*. ACM, 2018-08-27, pp. 1–6. [Online]. Available: <https://dl.acm.org/doi/10.1145/3230833.3230837>
- [18] M. S. Pour, E. Bou-Harb, Kavita Varma, N. Neshenko, D. A. Pados, and K.-K. R. Choo. "Comprehending the IoT cyber threat landscape: A data dimensionality reduction technique to infer and characterize internet-scale IoT probing campaigns," *Digital Investigation*, vol. 28, pp. S40–S49, 2019-04-01. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1742287619300246>
- [19] P. Richter and A. Berger. "Scanning the scanners: Sensing the internet from a massively distributed network telescope," in *Proceedings of the Internet Measurement Conference*, ser. IMC '19. Association for Computing Machinery, 2019-10-21, pp. 144–157. [Online]. Available: <https://doi.org/10.1145/3355369.3355595>
- [20] S. Kamoen. "Honeytrack: Persistent honeypot for the internet of things," 2018. [Online]. Available: <https://repository.tudelft.nl/islandora/object/uuid%3A344bd7aa-0a17-47dc-92fd-bd6f7e7b08c8>
- [21] C.-J. Wu, K. Yoshioka, and T. Matsumoto. "ThingGate: A gateway for managing traffic of bare-metal IoT honeypot," *Journal of Information Processing*, vol. 28, no. 0, pp. 481–492, 2020. [Online]. Available: https://www.jstage.jst.go.jp/article/ipsjip/28/0/28_481/_article
- [22] J. Franco, A. Aris, B. Canberk, and A. S. Uluagac. "A survey of honeypots and honeynets for internet of things, industrial internet of things, and cyber-physical systems," *arXiv:2108.02287 [cs]*, 2021-08-04. [Online]. Available: <http://arxiv.org/abs/2108.02287>
- [23] A. Vetterl and R. Clayton. "Bitter harvest: Systematically fingerprinting low- and medium-interaction honeypots at internet scale," in *12th USENIX Workshop on Offensive Technologies (WOOT 18)*. Baltimore, MD: USENIX Association, Aug. 2018. [Online]. Available: <https://www.usenix.org/conference/woot18/presentation/vetterl>
- [24] D. Moore. "Network telescopes: Tracking denial-of-service attacks and internet worms around the globe," in *Proceedings of the 17th Conference on Systems Administration (LISA 2003), San Diego, California, USA, October 26-31, 2003*, E. Frisch, Ed. USENIX, 2003.
- [25] M. Safaei Pour, A. Mangino, K. Friday, M. Rathbun, E. Bou-Harb, F. Iqbal, S. Samtani, J. Crichigno, and N. Ghani. "On data-driven curation, learning, and analysis for inferring evolving internet-of-things (IoT) botnets in the wild," *Computers & Security*, vol. 91, p. 101707, 2020-04-01. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167404819302445>
- [26] F. Shaikh, E. Bou-Harb, N. Neshenko, A. P. Wright, and N. Ghani. "Internet of malicious things: Correlating active and passive measurements for inferring and characterizing internet-scale unsolicited IoT devices," *IEEE Communications Magazine*, vol. 56, no. 9, pp. 170–177, 2018-09. [Online]. Available: <https://ieeexplore.ieee.org/document/8466375/>
- [27] A. Mangino, M. S. Pour, and E. Bou-Harb. "Internet-scale insecurity of consumer internet of things: An empirical measurements perspective," *ACM Transactions on Management Information Systems*, vol. 11, no. 4, pp. 21:1–21:24, 2020-10-12. [Online]. Available: <https://doi.org/10.1145/3394504>
- [28] AT&T. Alienvault open source siem (ossim). [Online]. Available: <https://cybersecurity.att.com/products/ossim>
- [29] Censys. [Online]. Available: <https://censys.io>
- [30] Shodan. [Online]. Available: <https://www.shodan.io>
- [31] H. Griffioen, T. M. Booi, and C. Doerr. "Quality evaluation of cyber threat intelligence feeds," in *Applied Cryptography and Network Security - 18th International Conference, ACNS 2020, Rome, Italy, October 19-22, 2020, Proceedings, Part II*, ser. Lecture Notes in Computer Science, M. Conti, J. Zhou, E. Casalicchio, and A. Spognardi, Eds., vol. 12147. Springer, 2020, pp. 277–296. [Online]. Available: https://doi.org/10.1007/978-3-030-57878-7_14
- [32] X. Bouwman, H. Griffioen, J. Egbers, C. Doerr, B. Klievink, and M. van Eeten. "A different cup of TI? the added value of commercial threat intelligence," in *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, Aug. 2020, pp. 433–450. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity20/presentation/bouwman>
- [33] Internet-Storm-Center. Dshield honeypot. [Online]. Available: <https://isc.sans.edu/honeypot.html>
- [34] B. Koshnaw and S. Furnell. "Assessing cyber security consumer support from technology retailers," *Computer Fraud & Security*, vol. 2022, no. 3, 2022.
- [35] IETF. Internet engineering task force. [Online]. Available: <https://www.ietf.org/>
- [36] E. Lear, D. Romascanu, and R. Droms. (2019-03) IETF RFC 8520 - manufacturer usage description (MUD) specification. [Online]. Available: <https://tools.ietf.org/html/rfc8520>
- [37] A. Hamza, D. Ranathunga, H. H. Gharakheili, M. Roughan, and V. Sivaraman. "Clear as MUD: Generating, validating and applying IoT behavioral profiles," in *Proceedings of the 2018 Workshop on IoT Security and Privacy*, ser. IoT S&P '18. Association for Computing Machinery, 2018-08-07, pp. 8–14. [Online]. Available: <https://doi.org/10.1145/3229565.3229566>
- [38] A. Hamza, H. H. Gharakheili, and V. Sivaraman. "Combining MUD policies with SDN for IoT intrusion detection," in *Proceedings of the 2018 Workshop on IoT Security and Privacy*, ser. IoT S&P '18. Association for Computing Machinery, 2018-08-07, pp. 1–7. [Online]. Available: <https://doi.org/10.1145/3229565.3229571>

- [39] A. Hamza, H. H. Gharakheili, T. A. Benson, and V. Sivaraman, "Detecting volumetric attacks on IoT devices via SDN-based monitoring of MUD activity," in *Proceedings of the 2019 ACM Symposium on SDN Research*, ser. SOSR '19. Association for Computing Machinery, 2019-04-03, pp. 36–48. [Online]. Available: <https://doi.org/10.1145/3314148.3314352>
- [40] A. Hamza, D. Ranathunga, H. H. Gharakheili, T. Benson, M. Roughan, and V. Sivaraman, "Verifying and monitoring iots network behavior using MUD profiles," *CoRR*, vol. abs/1902.02484, 2019. [Online]. Available: <http://arxiv.org/abs/1902.02484>
- [41] Anonymous. Mudscoptool. [Online]. Available: <https://anonymous.4open.science/r/MUDscope-BD68/>
- [42] ——. Mudscoptool dataset. [Online]. Available: https://mega.nz/folder/1tdGHSRC#MKP_tRnNv1mycM3-Mm6UQ
- [43] H. Touqeer, S. Zaman, R. Amin, M. Hussain, F. Al-Turjman, and M. Bilal, "Smart home security: challenges, issues and solutions at different iot layers," *The Journal of Supercomputing*, vol. 77, no. 12, pp. 14 053–14 089, 2021.
- [44] MITRE. Att&ck framework. [Online]. Available: <https://attack.mitre.org/>
- [45] NIST. Mud related resources. [Online]. Available: <https://www.nccoe.nist.gov/mud-related-resources>
- [46] Open Networking Foundation. (2012-09-06) OpenFlow switch specification (1.3.1). [Online]. Available: <https://opennetworking.org/wp-content/uploads/2013/04/openflow-spec-v1.3.1.pdf>
- [47] CISCO. Software-defined networking. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/software-defined-networking/overview.html>
- [48] S. K. Tayyaba, M. A. Shah, O. A. Khan, and A. W. Ahmed, "Software defined network (SDN) based internet of things (IoT): A road ahead," in *Proceedings of the International Conference on Future Networks and Distributed Systems*. ACM, 2017-07-19, pp. 1–8. [Online]. Available: <https://dl.acm.org/doi/10.1145/3102304.3102319>
- [49] CISCO. Netflow version 9 flow-record format. [Online]. Available: https://www.cisco.com/en/US/technologies/tk648/tk362/technologies_white_paper09186a00800a3db9.html
- [50] H. J. Griffioen, "Scanners: Discovery of distributed slow scanners in telescope data," 2018. [Online]. Available: <https://repository.tudelft.nl/islandora/object/uuid%3Adcb1669d-d81e-4aa3-bbd1-65049c3209c5>
- [51] H. Griffioen and C. Doerr, "Examining mirai's battle over the internet of things," in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '20. Association for Computing Machinery, 2020-10-30, pp. 743–756. [Online]. Available: <https://doi.org/10.1145/3372297.3417277>
- [52] M. Antonakakis, T. April, M. Bailey, M. Bernhard, E. Bursztein, J. Cochran, Z. Durumeric, J. A. Halderman, L. Invernizzi, M. Kallitsis *et al.*, "Understanding the mirai botnet," in *26th USENIX security symposium (USENIX Security 17)*, 2017, pp. 1093–1110.
- [53] H. Kang, D. H. Ahn, G. M. Lee, J. D. Yoo, K. H. Park, and H. K. Kim, "Iot network intrusion dataset." IEEE Dataport, 2019. [Online]. Available: <https://dx.doi.org/10.21227/q70p-q449>
- [54] L. McInnes and S. Horn. (2017-11-10) How HDBSCAN works — hdbscan 0.8.1 documentation. [Online]. Available: https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html
- [55] L. McInnes and M. Suppa. (2018-12-31) Benchmarking performance and scaling of python clustering algorithms — hdbscan 0.8.1 documentation. [Online]. Available: https://hdbscan.readthedocs.io/en/latest/performance_and_scalability.html
- [56] S. Guo, W. Lin, K. Zhao, and Y. Su, "Comparison of Clustering-based Network Traffic Anomaly Detection Methods," in *2021 IEEE 4th Advanced Information Management, Communicates, Electronic and Automation Control Conference (IMCEC)*, vol. 4, Jun. 2021, pp. 365–369, iSSN: 2693-2776.
- [57] R. Pang, V. Yegneswaran, P. Barford, V. Paxson, and L. Peterson, "Characteristics of internet background radiation," in *Proceedings of the 4th ACM SIGCOMM conference on Internet measurement - IMC '04*. ACM Press, 2004, p. 27. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1028788.1028794>
- [58] Wireshark. [Online]. Available: <https://www.wireshark.org/>
- [59] tcpdump. [Online]. Available: <https://www.tcpdump.org/>
- [60] Scapy. [Online]. Available: <https://scapy.net/>
- [61] Ubuntu. nfdump manual. [Online]. Available: <https://manpages.ubuntu.com/manpages/xenial/man1/nfdump.1.html>
- [62] Scikit-learn. [Online]. Available: <https://scikit-learn.org/stable/>
- [63] Numpy. [Online]. Available: <https://numpy.org/>
- [64] Anonymous. Mudscoptool results set. [Online]. Available: <https://mega.nz/folder/hx8VgRxa#9tBD8Mh8DpllfzobQcF45w>
- [65] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance," *The Journal of Machine Learning Research*, vol. 11, pp. 2837–2854, 2010.
- [66] A. Rosenberg and J. Hirschberg, "V-measure: A conditional entropy-based external cluster evaluation measure," in *Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL)*, 2007, pp. 410–420.
- [67] D. Dodson, D. Montgomery, T. Polk, M. Ranganathan, M. Souppaya, S. Johnson, A. Kadam, C. Pratt, D. Thakore, M. Walker, E. Lear, B. Weis, W. C. Barker, D. Coclin, A. Hojjati, C. Wilson, T. Jones, A. Baykal, D. Cohen, K. Yeich, Y. Fashina, P. Grayeli, J. Harrington, J. Klosterman, B. Mulugeta, S. Symington, and J. Singh, "Securing small-business and home internet of things (IoT) devices: Mitigating network-based attacks using manufacturer usage description (MUD)," in *NIST SPECIAL PUBLICATION 1800-15*, 2021-05-25. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1800-15.pdf>
- [68] Y. M. P. Pa, S. Suzuki, K. Yoshioka, T. Matsumoto, T. Kasama, and C. Rossow, "Iotpot: Analysing the rise of iot compromises," in *9th USENIX Workshop on Offensive Technologies (WOOT 15)*. Washington, D.C.: USENIX Association, Aug. 2015. [Online]. Available: <https://www.usenix.org/conference/woot15/workshop-program/presentation/pa>
- [69] M. S. Pour, D. Watson, and E. Bou-Harb, "Sanitizing the IoT cyber security posture: An operational CTI feed backed up by internet measurements," in *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2021-06, pp. 497–506. [Online]. Available: <https://ieeexplore.ieee.org/document/9505129/>
- [70] C. J. T. M. Schutijser, "Towards automated DDoS abuse protection using MUD device profiles," 2018-08-30. [Online]. Available: <http://essay.utwente.nl/76207/>
- [71] S. N. Matheu, J. L. H. Ramos, S. Pérez, and A. F. Skarmeta, "Extending MUD profiles through an automated iot security testing methodology," *IEEE Access*, vol. 7, pp. 149 444–149 463, 2019. [Online]. Available: <https://doi.org/10.1109/ACCESS.2019.2947157>
- [72] S. N. Matheu, A. Molina Zarca, J. L. Hernández-Ramos, J. B. Bernabé, and A. S. Gómez, "Enforcing behavioral profiles through software-defined networks in the industrial internet of things," *Applied Sciences*, vol. 9, no. 21, p. 4576, 2019-10-28. [Online]. Available: <https://www.mdpi.com/2076-3417/9/21/4576>
- [73] S. M. Sajjad, M. Yousaf, H. Afzal, and M. R. Mufti, "eMUD: Enhanced manufacturer usage description for IoT botnets prevention on home WiFi routers," *IEEE Access*, vol. 8, pp. 164 200–164 213, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9187209/>
- [74] A. Feraudo, P. Yadav, V. Safronov, D. A. Popescu, R. Mortier, S. Wang, P. Bellavista, and J. Crowcroft, "CoLearn: enabling federated learning in MUD-compliant IoT edge networks," in *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, ser. EdgeSys '20. Association for Computing Machinery, 2020-04-27, pp. 25–30. [Online]. Available: <https://doi.org/10.1145/3378679.3394528>
- [75] Google. Definition of federated learning. [Online]. Available: <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>
- [76] Y. Afek, A. Bremner-Barr, D. Hay, R. Goldschmidt, L. Shafir, G. Avraham, and A. Shalev, "Nfv-based iot security for home networks using MUD," in *NOMS 2020 - IEEE/IFIP Network Operations and Management Symposium, Budapest, Hungary, April 20-24, 2020*. IEEE, 2020, pp. 1–9. [Online]. Available: <https://doi.org/10.1109/NOMS47738.2020.9110329>

A TRAFFIC FLOW FEATURES FROM NFDUMP

Selected non-protocol-specific features from the nfdump tool, at <https://manpages.ubuntu.com/manpages/xenial/man1/nfdump.1.html>.

- ts. Start Time - first seen;
- te. End Time - last seen;
- td. Duration;
- pr. Protocol;
- sa. Source address;
- da. Destination address;
- sp. Source port;
- dp. Destination port;
- sas. Source autonomous system;
- pas. Previous autonomous system;
- ipkt. Input packets;
- opkt. Output packets;
- ibyt. Input bytes;
- obyt. Output bytes;
- flg. TCP flags;
- dir. Direction: ingress, egress;
- bps. Bytes per second;
- pps. Packets per second;
- bpp. Bytes per packet.

B IOT INTRUSION DETECTION DATASET

By Kang et al. [53], overviewed in Table 11.

Traffic scenario	(→ = attacks)		# packets c.	
	Interested devices	Attack category	Total	Attack
benign	both	None	137k	-
dos-synflood	server→EZVIZ	SYN Flooding	106k	48k
dos-synflood	server→NUGU	SYN Flooding	35k	17k
scan-hostport	server→EZVIZ	Port Scanning	80k	5k
scan-hostport	server→NUGU	Port Scanning	19k	6k
scan-portos	server→EZVIZ	OS Detection	186k	4k
scan-portos	server→NUGU	OS Detection	24k	8k
mirai-udpflood	EZVIZ→server	UDP Flooding	592k	475k
mirai-udpflood	NUGU→server	UDP Flooding	592k	475k
mirai-ackflood	EZVIZ→server	ACK Flooding	156k	38k
mirai-ackflood	NUGU→server	ACK Flooding	156k	38k
mirai-httpflood	EZVIZ→server	HTTP Flooding	124k	5k
mirai-httpflood	NUGU→server	HTTP Flooding	124k	5k
mirai-bruteforce	EZVIZ→NUGU	Telnet Bruteforce	273k	1.5k
mirai-bruteforce	NUGU→EZVIZ	Telnet Bruteforce	180k	1k

Table 11: Description of the used portion of the IoT Network Intrusion Dataset, by Kang et al. [53]. Each traffic scenario is a traffic capture file containing related activities. The target device is either victim or executor (→) of the attack.

C PRELIMINARY MUD ANOMALY PREVENTION

Table 12 overviews the result of packet-filtering with MUD profiles, for the Kang dataset. Where the MRT packets are double those of the ground truth, this happens simply because the generated MUD profiles captured both incoming and outgoing packets, whereas in the dataset, only incoming packets are marked as ground truth. Additional packets regards probes from Microsoft, Amazon, Google

and alike. Overall, the results of the filtering set the base for our further analyses.

D CLUSTERING PERFORMANCE ON LESS REPRESENTED ATTACKS

We report below the outputs from clustering on the less represented attacks of mirai UDP flooding, and bruteforce cases.

In the first listing below, brute-force attacks have been left separated in *attacker* and *victim* scenarios.

```
{'MIRAI-HOSTBRUTEFORCE-ATTACKER': {
  'out-represented_avg_percentage': 33.52007469654529,
  'represented_avg_percentage': 62.582010582010575
},
'MIRAI-HOSTBRUTEFORCE-VICTIM': {
  'out-represented_avg_percentage': 35.63083566760038,
  'represented_avg_percentage': 76.19047619047619
},
'MIRAI-UDPFLOODING': {
  'out-represented_avg_percentage': 38.23529411764706,
  'represented_avg_percentage': 100.0
}
```

Clusters purity: 84.58937198067632

```
Clusters {label: [majority label percentage, majority label]}
{'-1': [0, ''],
'0': [100.0, 'MIRAI-UDPFLOODING'],
'1': [100.0, 'MIRAI-UDPFLOODING'],
'10': [60.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'11': [100.0, 'MIRAI-UDPFLOODING'],
'12': [100.0, 'MIRAI-UDPFLOODING'],
'13': [58.333333333333336, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'14': [100.0, 'MIRAI-UDPFLOODING'],
'15': [100.0, 'MIRAI-UDPFLOODING'],
'16': [100.0, 'MIRAI-UDPFLOODING'],
'17': [100.0, 'MIRAI-UDPFLOODING'],
'18': [100.0, 'MIRAI-UDPFLOODING'],
'19': [100.0, 'MIRAI-UDPFLOODING'],
'2': [100.0, 'MIRAI-UDPFLOODING'],
'20': [50.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'21': [64.28571428571429, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'22': [100.0, 'MIRAI-UDPFLOODING'],
'23': [100.0, 'MIRAI-UDPFLOODING'],
'24': [100.0, 'MIRAI-UDPFLOODING'],
'25': [100.0, 'MIRAI-UDPFLOODING'],
'26': [100.0, 'MIRAI-UDPFLOODING'],
'27': [100.0, 'MIRAI-UDPFLOODING'],
'28': [100.0, 'MIRAI-UDPFLOODING'],
'29': [100.0, 'MIRAI-UDPFLOODING'],
'3': [50.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'30': [100.0, 'MIRAI-UDPFLOODING'],
'31': [100.0, 'MIRAI-UDPFLOODING'],
'32': [100.0, 'MIRAI-UDPFLOODING'],
'33': [100.0, 'MIRAI-UDPFLOODING'],
'34': [100.0, 'MIRAI-UDPFLOODING'],
'35': [77.77777777777777, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'36': [80.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'37': [75.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'38': [85.71428571428571, 'MIRAI-HOSTBRUTEFORCE-VICTIM'],
'39': [60.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'4': [50.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'40': [66.66666666666667, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'41': [66.66666666666667, 'MIRAI-HOSTBRUTEFORCE-VICTIM'],
'42': [66.66666666666667, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'43': [80.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'44': [50.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'5': [50.0, 'MIRAI-HOSTBRUTEFORCE-ATTACKER'],
'6': [100.0, 'MIRAI-UDPFLOODING'],
'7': [100.0, 'MIRAI-UDPFLOODING'],
'8': [100.0, 'MIRAI-UDPFLOODING'],
'9': [100.0, 'MIRAI-UDPFLOODING']}
```

In this second listing, victim and attacker brute-force scenarios were merged into one brute-force label.

```
{'MIRAI-BRUTEFORCE': {
  'out-represented_avg_percentage': 61.76470588235294,
  'represented_avg_percentage': 100.0
},
'MIRAI-UDPFLOODING': {
  'out-represented_avg_percentage': 38.23529411764706,
```



```
[ 0 ] Device ID      Signature
      transitions window  Signature time window      MRT Feed
-----
11-eufy-doorbell [17 - 20]                2022-06-24 20:58:32
- 2022-06-24 20:59:52 clusters-evols-record-session4_11-eufy-doorbell
11-foscam-cam2 [7 - 10]                   2022-06-24 20:53:48
- 2022-06-24 20:54:38 clusters-evols-record-session4_11-foscam-cam2
Max features correlation : 1.0 --- Avg features correlation :
3.9832882741253074e-05 --- Combined score : 0.5000199164413707
Correlation penalty multiplier for clusters difference : 1
```

```
Correlation values for signature features:
{'all_dists_avg': -0.9996016711725875, 'mutual_matches_n':
 0.0, 'mutual_matches_percentage': 0.0, 'fwd_matches_n': 0.0,
'fwd_matches_percentage': 0.0, 'fwd_matches_agglomeration_avg':
0.0, 'bwd_matches_n': 0.0, 'bwd_matches_percentage': 0.0,
'bwd_matches_agglomeration_avg': 0.0, 'clusters_balance': 1.0}
```

```
[ 1 ] Device ID      Signature
      transitions window  Signature time window      MRT Feed
-----
11-foscam-cam [10 - 18]                   2022-06-24 20:54:59
- 2022-06-24 20:58:44 clusters-evols-record-session4_11-foscam-cam
11-foscam-cam2 [19 - 27]                  2022-06-24 20:59:38
- 2022-06-24 21:02:59 clusters-evols-record-session4_11-foscam-cam2
Max features
correlation : 0.9998584712478614 --- Avg features correlation
: 0.8650292909226683 --- Combined score : 0.9324438810852649
```

Correlation penalty multiplier for clusters difference : 1

```
Correlation values for signature features:
{'fwd_matches_agglomeration_avg': 0.5819416382552827,
'fwd_matches_n': 0.75333824386712, 'clusters_balance':
0.7718449849879597, 'bwd_matches_n': 0.8485552916276634,
'fwd_matches_percentage': 0.8862326436248709,
mutual_matches_n': 0.8978247615364664, 'mutual_matches_percentage':
0.9501808393990387, 'bwd_matches_percentage':
0.9615988290998247, 'bwd_matches_agglomeration_avg':
0.9989172055805962, 'all_dists_avg': 0.9998584712478614}
```

F AMI FORMULA

To produce the AMI score, we first compute the mutual information $MI(X, Y)$:

$$MI(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right)$$

For our analysis, x and y are the values of the flow features, and the attack class, respectively. We then adjust the value by normalizing it on the expected value for MI , and the average over the entropy of X and Y :

$$AMI(X, Y) = \frac{MI(X, Y) - E[MI(X, Y)]}{\max(H(X), H(Y)) - E[MI(X, Y)]}$$