# Poster: Shedding Light Into the Darknet: Scanning Characterization and Detection of Temporal Changes

Rupesh Prajapati, Vasant Honavar, Dinghao
Wu, John Yen
{rxp338,vuh14,dwu12,juy1}@psu.edu
Penn State University

Michalis Kallitsis
mgkallit@merit.edu
Merit Network, Inc.

## ABSTRACT

Network telescopes provide a unique window into Internet-wide malicious activities associated with malware propagation, denial of service attacks, network reconnaissance, and others. Analyses of this telescope data can highlight ongoing malicious events in the Internet which can be used to prevent or mitigate cyber-threats in real-time. However, large telescopes observe millions of events on a daily basis which renders the task of transforming this knowledge to meaningful insights challenging. In order to address this, we present a novel framework for characterizing Internet's background radiation and for tracking its temporal evolution. The proposed framework: (i) Extracts a high dimensional representation of telescope scanners composed of features distilled from telescope data and learns an information-preserving low-dimensional representation of these events that is amenable to clustering; (ii) Performs clustering of resulting representation space to characterize the scanners and (iii) Utilizes the clustering outcomes as "signatures" to detect temporal changes in the network telescope.

## 1 INTRODUCTION

Cyber-attacks present one of the most severe threats to the security of the nation's critical infrastructure. The ubiquitous nature of the Internet-of-Things has expanded the threat surface and the number of on-net devices that can be easily compromised. Hence, network *situational awareness* becomes a germane task for network operators. A critical phase in most cyber-attacks is "reconnaissance", which includes "scanning" for potentially vulnerable devices on the internet that can be exploited later. Early detection and accurate characterization of these scanning behaviors can reveal new malware, their propagation strategies, motives, emerging or existing vulnerabilities, and unauthorized use of Internet resources.

However, network situational awareness is a challenging task. Scanning activities are oftentimes low in volume and interleaved with dominant normal network traffic. Security practitioners can employ (distributed) honeypots [3] for detecting the onset of such activities. Nevertheless, employing a large cluster of honeypots can be expensive. Further, high-interaction honeypots could be hard to maintain and adapt. Alternatively, one can monitor traffic destined to a "Network Telescope". Network telescopes or "Darknets" receive and record unsolicited traffic destined to an unused but routed address space and hence, provide a unique opportunity for characterizing Internet-wide malicious activities in a *timely manner*. This "dark IP space" hosts no services or devices, and therefore, any traffic arriving to it is inherently malicious. Darknets have been frequently used by the networking and security communities to shed light into dubious malware propagation other types of attacks[2].

## 2 PROBLEM FORMULATION

In this work, we aim to shed more light into the rapidly changing Darknet ecosystem (e.g., for an illuminating trend, see [5], Fig. 13) to facilitate accurate and timely situational awareness (i.e., identify zero-day exploits and new attacks in real time). **Our primary objective is to identify the patterns of scanning behaviors and track dynamic changes in these patterns**.

Each source IP, observed in the Darknet on a given day, can be characterized using attributes extracted from the raw Darknet data and other data sources (e.g., metadata from Censys.io, geolocation data, etc.). Scanning patterns can be captured by performing **clustering** on these attributes such that source IPs associated to a particular internet event get grouped together. Applying distance-based clustering algorithms on a high-dimensional data like the one collected in our Darknet results in unstable clusters. This can be overcome by employing **dimensionality reduction** techniques which learn information-preserving low-dimensional representation of such wide data. The meaningful clusters obtained on this low-dimensional embedding can be utilized to study and understand ongoing scanning events. More importantly, the task of **detecting temporal changes in scanning behaviors** can now be cast as a *goodness-of-fit* problem that compares the cluster signatures/distributions between consecutive days.

## 3 METHODOLOGY

Our starting point is raw traffic traces from a large network telescope (our team has access to a /13 Darknet spanning approximately 500,000 unique IPs; on a typical day, more than 100 GB of compressed Darknet data is collected consisting of some 3 billion packets on average). The nature of the clustering outcome is largely determined by the feature set chosen to characterize

the scanners. For this project, we extract a set of relevant and interpretable features, such as those that describe the volume and duration of scanning, the destination host and subnet dispersion, the set of ports and protocols scanned, and the set of services open at the scanners themselves as discovered by `Censys.io`.

**Task 1: Autoencoders for dimensionality reduction.** The extracted feature space is very high dimensional (consider, for example, that there are $2^{16}$ unique ports that needs to be considered). Clustering of such high-dimensional data poses a challenge in itself because distance calculations are known to be inherently unreliable in high-dimensional settings [1], making it challenging to apply standard clustering methods that rely on measuring distance between data samples to cluster them. Motivated by the recent success of deep representation learning, we employ Multi-layer Perceptron (MLP) autoencoders to learn the low-dimensional numerical embeddings of the input data.

Assuming $N$ sources and a feature vector of size $P$, our objective is to learn a *low-dimensional representation* of size $Q << P$ for the data that can be used for the clustering task introduced next. More formally, our goal is to learn functions $f(\cdot) : \mathbb{R}^P \to \mathbb{R}^Q$ and $g(\cdot) : \mathbb{R}^Q \to \mathbb{R}^P$ that minimize the reconstruction loss: $\sum_{i=1}^{N} (\ell(g \circ f(\mathbf{x}_i), \mathbf{x}_i))$, where the $\ell(\cdot) : \mathbb{R}^P \to \mathbb{R}$ is a loss function.

**Task 2: Clustering via k-means.** In order to characterize the Darknet behavior, the low-dimensional embeddings $f(\mathbf{x}_i) := \mathbf{z}_i \in \mathbb{R}^Q$ learnt in Task 1 are then assigned into groups based on similarity metrics. This can be formulated as the classical *clustering problem* of assigning $N$ objects into $K$ clusters such that some loss function is minimized. There are several approaches that one can adopt here; for this work, we elected to work with *k-means*. Note that, the number of clusters ($K$) is not known *a priori* and one needs to apply data-driven heuristics to find an appropriate value.

**Task 3: Earth Mover's Distance for change-point detection.** We employ the Earth Mover's Distance [6] (EMD) metric to measure the dissimilarity between the clustering "signature" of day $t$ and day $t + 1$. In our setting, each clustering outcome defines a distribution or "signature" that can be used for comparisons. If we denote the set of clusters obtained from clustering for day $t$ as $\{C_{1t}, C_{2t}, \dots, C_{Kt}\}$ and the centers of all clusters as $\{m_{1t}, m_{2t}, \dots, m_{Kt}\}$ where $m_{it} = \frac{\sum_{j \in C_{it}} x_{jt}}{|C_{it}|}$, $i = 1, \dots, K$, and $x_{jt} \in \mathbb{R}^P$, $j = 1, \dots, N$. Then, the signature $S_t = \{(m_{1t}, w_{1t}), \dots, (m_{Kt}, w_{Kt})\}$ can be employed, where $w_i$ represents the "weight" of cluster $i$ which is equal to the fraction of items in that cluster over the total population of scanners.

In our work, we posit the following "goodness-of-fit" problem: $H_0 : S_t = S_{t+1}$ versus $H_1 : S_t \neq S_{t+1}$. We reject the null hypothesis (i.e., we detect a change-point) when the EMD distance between the two distributions $S_t$ and $S_{t+1}$ exceeds a user-defined threshold (i.e., it is above an empirically-found baseline). As part of our ongoing work, we seek to find more statistically rigorous methods for this hypothesis testing problem.

## 4 EXPOSITION OF EARLY RESULTS

We showcase our methodology in Fig. 1 where we applied the proposed framework during the month of September, 2020. The figure shows a significant increase in the EMD distance between the clustering outcome of Sep. $5^{th}$ and Sep. $6^{th}$ which indicates a structural change in our Darknet. This increase in EMD distance is
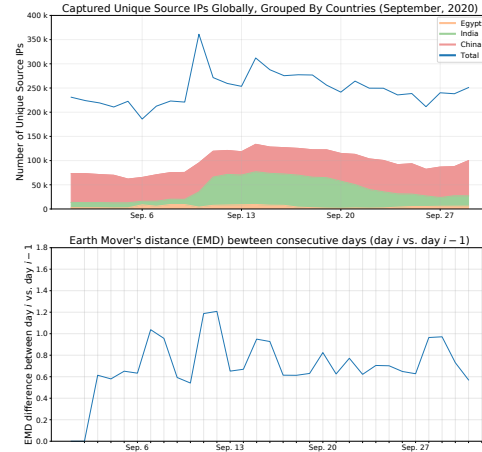


**Figure 1: Detecting temporal changes using the EMD.**

corroborated by the sudden increase in scanning traffic originating from countries like India and Egypt. In our clustering results for Sep. $6^{th}$, we observe novel clusters with scanners scanning a particular set of ports: 23, 80, 2323, 7574, 8080, 37215, 49152 and 52869. On cross-checking these particular port-scanning behavior with other sources (i.e., honeypot data and online reports), we were able to associate this traffic to the Mozi botnet. We also applied our techniques in the months of Nov. 2020 and Jan. 2021 and found other important events; for space economy, we omit the results.

## 5 CONCLUSIONS

We presented a novel framework towards network situational awareness. In addition to Darknet characterization (also done in other works, e.g., [4]), our approach utilizes the clustering outcomes to detect structural changes in the Darknet. Timely detection of such behavior would lead to rapid mitigation of emerging threats (e.g., zero-day exploits). As part of ongoing work, we plan to expand the set of features we select (e.g., introduce some of the ones in [3, 5]) to enhance the clustering interpretation. Moreover, given the limitations of running a centralized Darknet sensor [5], we plan to integrate additional data sources into our system (e.g., distributed honeypots, VirusTotal, ExploitDB, etc.) to further validate our results and to apply our techniques to other critical data sources.

## REFERENCES

[1] Charu C. Aggarwal, Alexander Hinneburg, and Daniel A. Keim. 2001. On the surprising behavior of distance metrics in high dimensional space. *International conference on database theory*, 420–434.

[2] Manos Antonakakis *et al.* 2017. Understanding the mirai botnet. In *26th USENIX Security Symposium (USENIX Security 17)*. 1093–1110.

[3] Paul Barford, Yan Chen, Anup Goyal, Zhichun Li, Vern Paxson, and Vinod Yegneswaran. 2010. *Employing Honeynets For Network Situational Awareness*. Springer US, Boston, MA, 71–102. https://doi.org/10.1007/978-1-4419-0140-8_5

[4] Félix Iglesias and Tanja Zseby. 2017. Pattern discovery in internet background radiation. *IEEE Transactions on Big Data* (2017).

[5] Philipp Richter and Arthur Berger. 2019. Scanning the Scanners: Sensing the Internet from a Massively Distributed Network Telescope. In *IMC'19*. 144–157.

[6] Yossi Rubner, Carlo Tomasi, and Leonidas J. Guibas. 1998. A Metric for Distributions with Applications to Image Databases. In *ICCV '98*.