

# Anomaly Detection in IoT Devices using Provenance Graphs

Ebelechukwu Nwafor

Howard University

Washington, DC 20059

ebelechukwu.nwafor@bison.howard.edu

Gedare Bloom

Howard University

Washington, DC 20059

gbloom@howard.edu

## ABSTRACT

The Internet of Things (IoT) has revolutionized the way we interact with devices. From smart grids, to healthcare and home automation systems. This paradigm shift inadvertently allows for ease of device access. Unfortunately, this technological advancement has been met with unforeseen security challenges. One major challenge is malicious intrusions –A common way of detecting a malicious attack is by treating an attack as an anomaly and using anomaly detection techniques to pinpoint the source of an intrusion. In a given IoT device, provenance graphs which denotes causality between system events offers immense benefit for anomaly detection. Provenance provides a comprehensive history of activities performed on a system which indirectly ensures trust. Given a provenance graph, how do we determine if anomalous activities exists? This paper seeks to address this issue. In this paper, We introduce an error tolerant graph embedding technique using frequencies of nodes and edges in which provenance graphs are converted into a vector space representation. This vector space representation of the graphs from learning and detection phase is used as an input parameter to our anomaly detection algorithm. we propose two anomaly detection algorithms using provenance graphs. The first approach involves the use of a similarity metric to compare provenance graphs while the later involves the use of a hybrid anomaly detection algorithm. We evaluate the effectiveness of our approach with IoT applications which generates provenance graphs.

## CCS CONCEPTS

•Computer systems organization →Embedded systems; *Internet of Things*; Data Provenance;

## KEYWORDS

IoT, Data Provenance, Anomaly Detection

### ACM Reference format:

Ebelechukwu Nwafor and Gedare Bloom. 2016. Anomaly Detection in IoT Devices using Provenance Graphs. In *Proceedings of ACM Conference, Washington, DC, USA, July 2017 (Conference'17)*, 6 pages. DOI: 10.475/123.4

## 1 INTRODUCTION

Over the years, there has been an increase in the number of intrusions reported on IoT devices. With the pervasive nature of these

smart devices, the complexity of intrusions are further exacerbated. Malicious intrusions on IoT devices could have disastrous financial consequences. For example, a vulnerability on a consumer device such as a smart watch or mobile card reader could lead to the theft of sensitive financial and personal information. One way of detecting intrusions is by the use of anomaly detection techniques. An anomaly, also referred to as an outlier, is defined as data that deviates from the normal system behavior. This enables the detection of known and unknown malicious attacks. Anomaly detection has applications in domains such as intrusion detection, fraud detection, medical health devices, sensor fault detection, web spam. An anomalous event could indicate that a system fault exists. It could also indicate that a system is being used as a botnet in a distributed denial of service attack (DDOS). Due to the sensitive nature of safety critical systems, detecting current and future malicious attacks is of utmost importance to the security of IoT devices. However, ensuring data trust is a challenging task. How do we provide an effective means of detecting anomalous instances in a given system? Provenance can be used to address this issue. Provenance graphs captures a holistic history of system events and also offers an efficient way of representing relationships between multiple data objects which can be used to detect system faults or anomalous system behaviors. For example, sensors deployed in a oil rig, provenance can be used to detect when there is a device leak. Also, in an IoT enabled smart home, provenance can be used to detect a point of intrusion in an event of a system hack.

In this paper, we motivate the need for device anomaly detection on memory constrained IoT devices. Unlike most anomaly detection system which utilize system call frequencies to detect anomalous system behaviors, We take a different approach to anomaly detection by detecting anomalous system events in IoT devices based on provenance data generated by these devices. We identify how provenance graphs can be used to detect anomalous data instances. We provide an overview of anomaly detection, and provenance graphs. We propose a lightweight graph similarity algorithm for detection of anomalous data instances in IoT devices. This algorithm relies on measuring the similarity of provenance graphs in the learning and detection phase. Our technical contributions in this paper are outlined in detail as follows:

- We introduce a graph-vector space approach to representing provenance graphs in vector space. This allows the comparison of graphs using similarity metrics such as cosine similarity or jaccard distance.
- We propose two anomaly detection techniques for IoT devices using provenance graphs. The first algorithm is based on the similarity between provenance graphs contained in both the learning and detection phase which is achieved by using a similarity metric. A threshold is set which classifies a provenance graph as anomalous once the threshold is exceeded. The second

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

Conference'17, Washington, DC, USA

© 2016 Copyright held by the owner/author(s). 123-4567-24-567/08/06...\$15.00  
DOI: 10.475/123.4

algorithm is a hybrid approach which involves the uses of clustering algorithm, DBSCAN, and  $k$ -nearest neighbor algorithm,  $k$ -nearest neighbor.

- Finally, we evaluate the effectiveness of our approach to graph based anomaly detection with data set from sample IoT application(s). Detailed results are presented in sections IV.

The remaining portion of this paper is organized as follows. In Section 2, we discuss background information on Anomaly detection, Provenance Graphs, text based query ranking and a review of  $k$ -NN. In Section 3, we discuss how provenance data is converted in vector space using tf-idf. Section 5 talks about our classification approach using  $k$ -NN. Section 6 describes our experiment and gives a detailed evaluation. We summarize and conclude in Section 7 and Section 8 contains further discussions.

## 2 PROBLEM STATEMENT

Given a set of provenance graphs,  $P = \{p_1, p_2, \dots, p_n\}$  containing provenance data from both the learning and detection phase, we want detect anomalies that might exist in  $P$ . This involves converting elements in  $P$  into a feature vector space. Anomalies are detected by comparing the similarity of graphs contained in  $P$  using a similarity metric and also by using a hybrid approach as discussed in section IV

## 3 BACKGROUND

This section describes key concepts of Anomaly detection techniques, Provenance graphs, text-based query ranking and  $k$ -NN.

### 3.1 Overview of Anomaly Detection

The definition of an anomaly is often domain specific. Hawkins, a primer anomaly detection researcher defines an anomaly as an "observation which deviates so much from the other observations as to arouse that it was generated by a different mechanism [4]. An Anomaly is usually considered to be in isolation from similar data points. The notion of a normal behavior is defined by the average behavior that is depicted where most data points are centered in. Anomaly detection has been researched in a wide variety of fields such as statistics, machine learning, information theory. It has applications in finance for fraud detection, in health care for the monitoring of patient care, in computer security for intrusion detection. For the purpose of this paper, we focus on anomaly detection on memory constrained IoT devices. An anomaly often indicates the presence of a malicious entity or a system fault. For example, an anomaly could be a sudden increase in web traffic of a web server. This could be indicative of a denial of service attack. Additionally, In a critical care health device such as a pacemaker, an anomalous event could be detrimental to the health of a patient which could result in the loss of lives.

Most of the work done on anomaly detection is centered on the detection of anomalous behavior in point based datasets. This might ignore the dependent relationship that might exist between data points. Graphs provides a means of modeling complex structures such as social networks, computer networks, DNA sequences.

The process of determining all anomalous instances in a given dataset or system is a complex task. The challenge in anomaly detection is providing the right feature from a dataset to use for

detection. Another challenge exists in defining what constitutes as normal system behavior. There often exist a thin line between what is considered normal system behavior and what is considered an anomaly. In addition, what is considered normal system behavior is constantly evolving. The issue of generating training or test dataset which classifies anomalous and normal system behavior is a major challenge since not all known anomalous system behavior can be recreated.

Anomaly detection consists of two phase: Learning or training phase and the test or detection phase. In the training phase, the system collects training dataset. This data is considered to be a representation of the system's normal daily activity and free from malicious events. Once training dataset has been collected, the system's activities are further observed. This part is known as the testing phase. In the testing phase, observed system behavior is compared to the Learning phase to determine if an anomaly exists between the two. A threshold as defined by domain experts is used to determine if the observed data is considered an anomaly.

Anomaly detection involves the use of statistical or machine learning techniques such as clustering and classification to determine normal or anomalous data instances. Some methods deal with assigning a score to determine the anomaly. Details on anomaly detection techniques are outlined below

- Statistical-based approach: This approach involves the use of parametric and non parametric statistical inference to develop models which are used to determine if a dataset fits a statistical model. Instances that do not fit the defined statistical model are classified as an anomaly. Parametric method assumes the knowledge of a distribution. On the other hand, non parametric distribution does not assume prior distribution knowledge.
- Classification-based approach: The main idea in this approach involves building models which use training data set with pre-defined labels (i.e normal, anomalous) to classify incoming data. Classification works in two phase: training phase and observation phase. In the training phase, data is collected which contains labels of normal and anomalous system behavior. If the dataset only contains a label of either anomalous or normal behavior, this is considered as a one class classifier. In the observation phase, incoming data is classified by defined data labels.
- Clustering-based approach: The main idea is to group similar data instances into clusters. There are various approach to clustering. One approach looks at the density of the clusters, normal data belongs to large dense clusters while abnormal data belongs to small clusters. Another approach as treats clustering as one class which assumes that normal belongs to a cluster and abnormal data does not. Another approach looks at the distance of the data to the centroid. Centroids are seen as the center of the cluster. Normal data are considered to be closer to the centroid than anomalies lie
- Nearest-Neighbor based approach: It is based on the assumption that normal data occurs in dense neighborhoods and abnormal data in sparse neighborhoods. The main idea is to assign an incoming observation data to a class based on its proximity to the closest data point in the training data set. A distance or similarity measure is used to quantify the distance between points in a dataset. A popular form of nearest neighbor technique

is the  $k$ -nearest neighbor which groups incoming data based on proximity to  $k$  closest data point. Details on  $k$ -nearest neighbors is discussed in section 3.5.

- **Density-based approach:** This approach is used to estimate the density of  $k$  nearest neighbors. A data instance in a dense neighborhood is considered normal while data instances in neighborhoods with a sparse density are considered anomalous. The distance from a data instance to a nearest neighbor is seen as the inverse of the density of data instances. This approach faces an issue in which the density approach performs poorly in regions of varying densities. Local Outlier Factor (LOF) addresses this issue. LOF is a measure of the degree of Outlieriness of each data instance contained in a data set. It is achieved by comparing the ratio of local density of  $k$  nearest neighbors to the density of a data instance. Data instances with lower density are considered outliers.

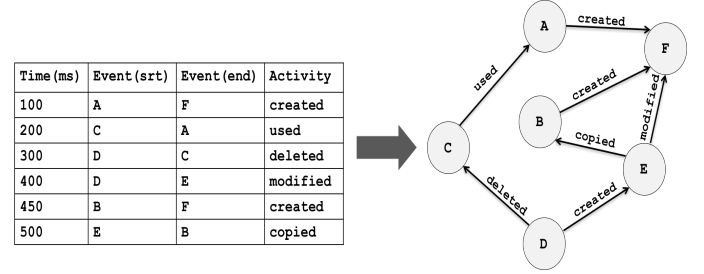
Data labels are grouped into three major categories. Supervised anomaly detection, semi-supervised anomaly detection, and unsupervised anomaly detection. In supervised anomaly detection approach, training data contains instances of normal and anomalous data. Incoming data is classified based on the training data category. In semi-supervised approach, only one class of training data is collected, normal data. Incoming data that does not fit the normal class as specified by a threshold is regarded as an anomalous. In the training phase, most anomaly detection techniques use data derived from the normal system behavior. This is referred to as one class classification. In unsupervised approach, there are no training dataset. It is believed that normal data is clustered around each and occurs more frequently than an anomaly. This is a widely used form of anomaly detection since training data which is hard to get is not required. Detailed information on anomaly detection techniques can be seen in [1, 3, 5, 16]

### 3.2 Provenance Graphs

Provenance denotes the origin of an object and all other activities that occurred. An example of provenance can be seen with a college transcript. A transcript is the provenance of a college degree because it outlines all of the courses satisfied in order to attain the degree. In the field of computing, data provenance, also known as data lineage, can be defined as the history of all activities performed on a data object from its creation to its current state. Provenance ensures trust of data [2]. It outlines causality and dependency between all objects involved in the system and allows for the verification of a data source. In an IoT device, a provenance graph  $G = (V, E)$  is a directed acyclic graph (DAG) in which vertices represent device or sensor events data and the edges correspond to the interaction between them.

Graphs provide a means of modeling complex relationships that exists between data objects and is a good choice for modeling provenance because provenance data contains information with dependency relationships which denotes causality between multiple data objects. Provenance data is represented by directed acyclic graphs where nodes denote data objects and the edges represent relationships between data objects. For example, provenance graph from a device can be generated by evaluating the log of system calls. This log information might contain noisy data and is further streamed

by mapping it to a provenance model. With this information, we are able to build a workflow of device data execution. There have been numerous provenance collection systems developed to track provenance in a computing device most of which deal with tracking system calls [9, 12, 15]. Provenance graph used for experimentation is generated from PAIoT, a provenance-aware framework [1]. We chose PAIoT because it captures dependencies between device and sensor events.



**Figure 1: Provenance data transcribed to Provenance graph which depicts causal dependency between system events. The nodes represent events while the edges represent activities**

Causality and dependency are concepts used to denote relationship between system events. Provenance graphs in turn can be used in digital forensics [16] to investigate the cause of a malicious attack and also in intrusion detection systems to further enhance the security of computing devices. For further reading on provenance models and provenance graphs, we refer the reader to [1].

### 3.3 Similarity Measures

Similarity defines how identical two objects are. It measures some form of distance between data objects by either measuring the angle between the objects (Cosine similarity) or a linear distance between the objects (Euclidean distance). Similarity measures are widely used in document retrieval for selecting a query given a list of documents. The similarity measure used is application dependent on the data set and the sparsity of data. Based on literature, there are three well known similarity measures for evaluating vectors. These measures are outlined below:

**3.3.1 Cosine similarity:** This is a measure of orientation between the two non-zero vectors. It measures the cosine of the angle between the vectors. Two vectors which are at an angle of  $90^\circ$  have a similarity of 0 while two vectors which are similar (with an angle of  $0^\circ$ ) have a cosine of 1 and two vectors which are completely opposite (with an angle of  $180^\circ$ ) have a similarity of -1. Since we are concerned with the similarity of the vectors, we are only concerned with the positive values of bounded in  $[0, 1]$ . To compute the cosine similarity between two vectors,  $X$  and  $Y$ , cosine similarity is represented by using the dot product and magnitude of the two vectors.

$$\cos(\theta) = \frac{X \cdot Y}{\|X\| \cdot \|Y\|} = \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n X_i^2} \sqrt{\sum_{i=1}^n Y_i^2}}$$

**3.3.2 Jaccard Similarity:** This similarity measure evaluates the intersection divided by the union of two non zero vectors.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

**3.3.3 Euclidean distance:** This measures calculates the line distance between two data objects in an euclidean space. The euclidean distance between vectors  $A$  and  $B$ ,  $d(A, B)$  is defined by:

$$d(A, B) = \sqrt{\sum_{i=1}^n (A_i - B_i)^2}$$

### 3.4 A Review of DBSCAN and $k$ -nearest neighbors

**3.4.1  $k$ -Nearest Neighbors:**  $k$ -NN is an instance based supervised learning algorithm for classifying data based on nearest neighbors. Data is grouped on its similarity to nearest neighbors where  $k$  denotes the number of neighbors in which the input data is compared to. A similarity measure such as euclidean distance, jaccard similarity, or cosine similarity is used to measure the distance between vectors.  $k$ -NN can be applied to classification and regression problems.

Researchers have proposed various modifications to the  $k$ -NN algorithm. Ramaswamy et al. [13] proposed a  $k$ -NN modification which calculates the sparseness estimates for vectors in a dataset. The vectors are sorted in increasing order according the distance from its  $k^{th}$  neighbor. Bay and Swabacher demonstrates how pruning irrelevant datapoint which are not considered anomalous can result in a linear complexity of nearest neighbor search for a randomized data. A threshold score is assigned which is based on the score of the weakest anomaly found. Pruning can be achieved by using the relative density of data points, an anomaly is believed to occur in a group of data points with low density.

**3.4.2 Density Based Spatial Clustering of Applications with Noise (DBSCAN):** DBSCAN is a density-based clustering algorithm that differentiates regions with high density from regions with low density. It defines two parameters  $\epsilon$ , and **MinPts**.  $\epsilon$  defines the maximum distance between two neighboring points in which they are considered to be in the same cluster. **MinPts** is the minimum number of points that can be contained in a cluster.

Let  $S = \{s_1, s_2, \dots, s_n\}$  be a set of point to be clustered.  $S$  consists of three point categories, a core point,  $p$ , a border point,  $b$  and a noise point,  $n$ . A core point is a point in which its **minPts** are within distance  $\epsilon$ . These are at the interior of the cluster. Border points are points on the edge of the cluster. A noise point (outlier) is any point that is neither a core point or a border point.

A point,  $l$  is directly reachable from a core point,  $p$  if there exist a path which is within distance of  $\epsilon$  from the core another point. (i.e  $\exists \{p_1, p_2, \dots, p_n\}$ , where  $|p_{i+1} - p_i| \leq \epsilon$ ).

A core point,  $p$  that is within distance  $\epsilon$  (i.e  $p \leq \epsilon$ ) is considered part of the same cluster. A border point is also considered part of a cluster if it close to a core point.

## 4 GRAPH BASED ANOMALY DETECTION

### 4.1 Potential Anomalies

Due to the ubiquitous nature of IoT devices, there are a wide array of potential vulnerabilities associated with them. We focus on select vulnerabilities of inconsistent sensor output. Inconsistency are analogous to spikes in sensor readings, constant data values, or faulty sensors, inconsistent device attributes.

### 4.2 Graph to Vector Space Conversion

**4.2.1 Feature Extraction.** In order to apply clustering algorithms or similarity measure to provide anomaly detection, we compute a vector space representation of our provenance graphs. This representation is used as input parameters to our hybrid anomaly detection algorithm. Selecting features from graph properties is an important task because we need to select features that preserves the order and details of each node and edges contained in the graph. We focus on the frequency of the nodes and edges contained in both graphs. Selecting the right features is important in ensuring optimal performance of our anomaly detection algorithms. Our approach not only utilizes the frequency of nodes but also consists of a damping factor which regulates the weight of nodes or edges that occur frequently and increases the weight of nodes and edges that occurs less frequently.

Given a set of provenance graphs  $P = \{p_1, \dots, p_n\}$ , where  $p_x = (V, E)$ .  $P$  consists of graphs both in the learning and detection phase. We denote the occurrence of edges and nodes contained in set  $P$ . Each graph in  $P$  represents a vector by using our graph embedding approach which draws emphasis from document query retrieval. This approach preserves the order of edges and allows updates of edges of incoming graphs. We formally define our approach as follows:

**Definition 4.1.** Let  $P = \{p_1, \dots, p_n\}$  where  $p_i = (V, E)$ , the vectorial representation of  $p_i$ ,  $v_i$  is the number of times each vertice,  $V_i$  and edges,  $E_i$  contained in  $P$  appears in the graph,  $p_i$ .

$$v_x = (freq(E_i, p_i), freq(V_i, p_i))$$

where  $freq$  denotes the occurrence of  $E_i$  in graph  $p_i$

The order of Edges and vertices can be found by taking a breadth first search transversal of the graphs.

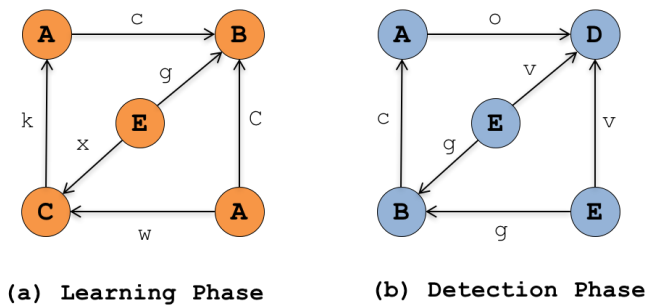


Figure 2: Two graphs one generated from the Learning phase and the other from the Detection phase

For example, Figure 2 displays provenance graph generated in the Learning and detection phase respectively. Both graphs,  $p_1$ , and  $p_2$  consists of vertices  $A, B, C, D$ , and  $E$ . The vector representation of the two graphs,  $v_1$ , and  $v_2$  are

$$v_1 = (1, 1, 1, 1, 0, 0)$$

$$v_2 = (1, 2, 0, 1, 1, 0)$$

### 4.3 Graph Based Similarity Detection Algorithm

The method of comparing the similarity of graphs based on a similarity metric is inspired by a document retrieval technique. Given a corpus  $D = \{d_1, \dots, d_n\}$ , and query,  $q$ . How do we find document(s)  $d_x, \dots, d_y$  which are similar to  $q$  and rank them by order of importance. To achieve this, documents are converted into a vector space representation which allows document to be ranked based on some similarity metric. Figure 3 depicts the overall goal of the similarity approach in detecting anomalies.

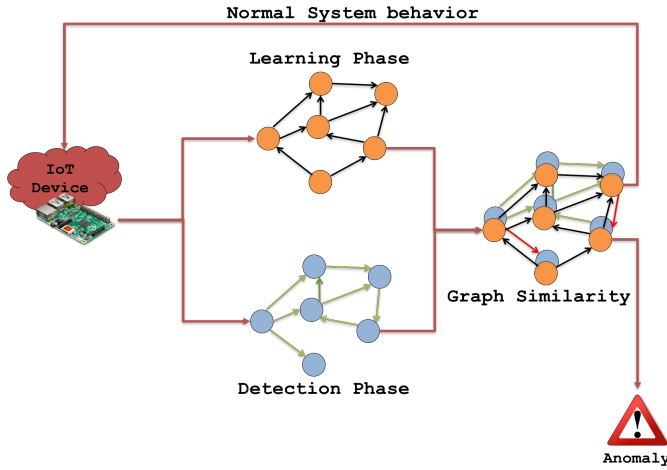


Figure 3: Graph Similarity Approach

Given  $v_x, v_y$  which denotes the vector representation of provenance graphs  $p_x, p_y$ . The similarity of  $p_x, p_y$  is found by calculating the cosine similarity between the two vectors where 1 denotes similarity between the two vectors and 0 denotes non-similarity between the two vectors. A threshold value is set which is used to classify the behavior of provenance graphs in the detection phase as normal or an anomaly.

$$\text{sim}(v_x, v_y) = \frac{v_x \cdot v_y}{\|v_x\| \cdot \|v_y\|} = \frac{\sum_{i=1}^n v_{xi} v_{yi}}{\sqrt{\sum_{i=1}^n v_{xi}^2} \sqrt{\sum_{i=1}^n v_{yi}^2}} \in [0, 1]$$

**4.3.1 Defining Anomaly Threshold:** An anomaly threshold  $t$  is a score that defines at what point a provenance graph contained in the test data is considered anomalous. Ensuring a proper threshold score is used for detection is an important task that requires extensive knowledge of the attack domain. The threshold is manually set to a value  $t$ , which is defined by domain experts. For automatic

anomaly threshold detection, one can use forecasting methods to define the anomaly score. Forecasting techniques are beyond the scope of this research.

4.3.2 Computational Complexity: **TODO**

### 4.4 Hybrid Detection Algorithm

**TODO...**

## 5 EXPERIMENTAL EVALUATION

**TODO...**

## 6 RELATED WORK

There has been a considerable amount of research done on anomaly detection. Most of the work which involves the analysis of system call events. Liao et al [7] characterizes a system's normal behavior by denoting the frequency of unique system calls which are converted to a vector space using the text classification approach. A classification algorithm  $k$ -NN is used to classify test data set. Our approach also deals with system events by converting a provenance graph to a vector representation. We also use a more sophisticated clustering and machine learning algorithm for the detection stage.

Some graph approach involves the use of a community based approach in which dense regions of connected nodes are considered normal and nodes with high sparse regions which do not belong to any community are considered anomalous. *AUTOPART* consist of nodes with similar neighbors are clustered together and the edges which do not belong to any cluster is considered as an anomaly. To find communities the graph achieves this task by reorganizing the rows and the columns of the adjacency list

Our approach looks at anomaly detection on graphs.

Additionally, anomaly detection on graphs has also been explored. Manzoor et al [8] proposed a centroid based clustering anomaly detection for instances of streaming heterogeneous graphs in real time. This algorithm is able to accommodate incoming edges in real time. They propose a method of comparing similarity between heterogeneous graphs by comparing the similarity of two graphs by their relative frequency. Each graphs is represented as a vector known as shingles. Since all of the graph is stored in memory, they also accommodate an efficient representation of the shingles in memory in what is known as streamhash.

Stephanie forest group [6] provided an analogy between the human immune system and intrusion detection which builds a normal system behavior repository by looking at the system call sequences of an application. This sequence is stores in a normal database which is queried for all other online behaviors are judged. If an application executes sequence of system calls that are not found in the normal database, a mismatch is recorded. If the mismatch for that application exceeds a threshold, an anomaly is detected.

Yoon [14] developed a technique for intrusion detection in embedded systems by analyzing system call frequencies. This is achieved by learning normal system profile by observing patterns in system call frequency distribution. Data from the training set is clustered using  $k$ -means to categorize legitimate system behavior. It is their belief that applications follow a known frequency pattern which

is centered around the centroid. Observation at run-time are compared with the clusters in the detection phase, if the incoming observation does not fit into a cluster, it is considered an anomaly.

Our approach is similar to graph kernels, and graph isomorphism. graph kernels involves measuring the similarity between two graphs, graph edit distance looks at the number of operations required for a graph  $G_1$  to be identical to  $G_2$ . graph isomorphism is a direct mapping of two graphs. More formally, two graphs are said to be isomorphic, if there exist a one to one mapping between edges and vertices contained in the graph.

Papadimitriou et al [11] proposed five similarity algorithms namely Signature similarity, vertex/edge vector similarity, vertex ranking, vertex edge overlap, for comparing the similarity of web graphs for the detection of anomalies inspired by document similarity method namely shingling and random projection based method. An anomaly could be a missing link (edge) or a web node. Nodes represents a webpage. Out of the five similarity measures proposed, Signature similarity which compares two graphs based on a set of features (signatures) using a scheme known as *simHash* performed the best followed by vector similarity. By comparing instances between snapshot of crawled webpages, this enables to detect inconsistencies in crawled we content.

Noble et al [10] proposed two algorithms for comparing graph similarity. The first approach, anomalous substructure detection looks at unusual substructures in graphs. This is achieved by inverting the measure of patterns that occurs frequently in a graph. The second approach examines

## 7 SUMMARY AND CONCLUSION

In this paper, we proposed two anomaly detection algorithms for intrusion detection on IoT devices. Our approach is lightweight and efficient in detecting anomalies which might exist using provenance graphs. We evaluated the functionality of our approach through implementation.

## 8 DISCUSSION

TODO

## 9 ACKNOWLEDGMENT

This research has been supported in part by US National Science Foundation (CNS grant No. 1646317), and by Leidos. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of NSF or Leidos.

## REFERENCES

- [1] Leman Akoglu, Hanghang Tong, and Danai Koutra. 2015. Graph based anomaly detection and description: a survey. *Data Mining and Knowledge Discovery* 29, 3 (01 May 2015), 626–688. DOI : <http://dx.doi.org/10.1007/s10618-014-0365-y>
- [2] Elisa Bertino. 2015. *Data Trustworthiness—Approaches and Research Challenges*. Springer International Publishing, Cham, 17–25. DOI : [http://dx.doi.org/10.1007/978-3-319-17016-9\\_2](http://dx.doi.org/10.1007/978-3-319-17016-9_2)
- [3] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly Detection: A Survey. *ACM Comput. Surv.* 41, 3 (July 2009), 15:1–15:58. DOI : <http://dx.doi.org/10.1145/1541880.1541882>
- [4] D. M. Hawkins. 1980. *Identification of outliers*. Chapman and Hall, London [u.a.]. [http://gso.gbv.de/DB=2.1/CMD?ACT=SRCHA&SRT=YOP&IKT=1016&TRM=ppn+02435757X&sourceid=fbw\\_bibsonomy](http://gso.gbv.de/DB=2.1/CMD?ACT=SRCHA&SRT=YOP&IKT=1016&TRM=ppn+02435757X&sourceid=fbw_bibsonomy)
- [5] Victoria Hodge and Jim Austin. 2004. A Survey of Outlier Detection Methodologies. *Artificial Intelligence Review* 22, 2 (Oct. 2004), 85–126. DOI : <http://dx.doi.org/10.1023/B:AIRE.0000045502.10941.a9>
- [6] Steven A. Hofmeyr, Stephanie Forrest, and Anil Somayaji. 1998. Intrusion Detection Using Sequences of System Calls. *J. Comput. Secur.* 6, 3 (Aug. 1998), 151–180. <http://dl.acm.org/citation.cfm?id=1298081.1298084>
- [7] Yihua Liao and V. Rao Vemuri. 2002. Using Text Categorization Techniques for Intrusion Detection. In *Proceedings of the 11th USENIX Security Symposium*. USENIX Association, Berkeley, CA, USA, 51–59. <http://dl.acm.org/citation.cfm?id=647253.720290>
- [8] Emaad Manzoor, Sadegh M. Milajerdi, and Leman Akoglu. 2016. Fast Memory-efficient Anomaly Detection in Streaming Heterogeneous Graphs. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*. ACM, New York, NY, USA, 1035–1044. DOI : <http://dx.doi.org/10.1145/2939672.2939783>
- [9] Kiran-Kumar Muniswamy-Reddy, David A. Holland, Uri Braun, and Margo Seltzer. 2006. Provenance-aware Storage Systems. In *Proceedings of the Annual Conference on USENIX '06 Annual Technical Conference (ATEC '06)*. USENIX Association, Berkeley, CA, USA, 4–4. <http://dl.acm.org/citation.cfm?id=1267359.1267363>
- [10] Caleb C. Noble and Diane J. Cook. 2003. Graph-based Anomaly Detection. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '03)*. ACM, New York, NY, USA, 631–636. DOI : <http://dx.doi.org/10.1145/956750.956831>
- [11] Panagiotis Papadimitriou, Ali Dasdan, and Hector Garcia-Molina. 2010. Web graph similarity for anomaly detection. *Journal of Internet Services and Applications* 1, 1 (01 May 2010), 19–30. DOI : <http://dx.doi.org/10.1007/s13174-010-0003-x>
- [12] Thomas Pasquier, Xueyuan Han, Mark Goldstein, Thomas Moyer, David Eyers, Margo Seltzer, and Jean Bacon. 2017. Practical Whole-System Provenance Capture. In *Symposium on Cloud Computing (SoCC'17)*. ACM, ACM.
- [13] Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. 2000. Efficient Algorithms for Mining Outliers from Large Data Sets. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data (SIGMOD '00)*. ACM, New York, NY, USA, 427–438. DOI : <http://dx.doi.org/10.1145/342009.335437>
- [14] Man-Ki Yoon, Sibin Mohan, Jaesik Choi, Mihai Christodorescu, and Lui Sha. 2017. Learning Execution Contexts from System Call Distribution for Anomaly Detection in Smart Embedded System. In *Proceedings of the Second International Conference on Internet-of-Things Design and Implementation (IoTDI '17)*. ACM, New York, NY, USA, 191–196. DOI : <http://dx.doi.org/10.1145/3054977.3054999>
- [15] Robert H'obbes' Zakon (Ed.). 2012. *28th Annual Computer Security Applications Conference, ACSAC 2012, Orlando, FL, USA, 3-7 December 2012*. ACM. <http://dl.acm.org/citation.cfm?id=2420950>
- [16] Y. Zhang, N. Meratnia, and P. Havinga. 2010. Outlier Detection Techniques for Wireless Sensor Networks: A Survey. *IEEE Communications Surveys Tutorials* 12, 2 (2010), 159–170. DOI : <http://dx.doi.org/10.1109/SURV.2010.021510.00088>