



# Graph clustering and anomaly detection of access control log for forensic purposes



Hudan Studiawan <sup>a, b, \*</sup>, Christian Payne <sup>a</sup>, Ferdous Sohel <sup>a</sup>

<sup>a</sup> School of Engineering and Information Technology, Murdoch University, Australia

<sup>b</sup> Department of Informatics, Institut Teknologi Sepuluh Nopember, Indonesia

## ARTICLE INFO

### Article history:

Received 15 July 2016

Received in revised form

21 March 2017

Accepted 1 May 2017

Available online 3 May 2017

### Keywords:

Authentication log

Improved MajorClust

Event log forensics

Anomaly detection

## ABSTRACT

Attacks on operating system access control have become a significant and increasingly common problem. This type of security threat is recorded in a forensic artifact such as an authentication log. Forensic investigators will generally examine the log to analyze such incidents. An anomaly is highly correlated to an attacker's attempts to compromise the system. In this paper, we propose a novel method to automatically detect an anomaly in the access control log of an operating system. The logs will be first preprocessed and then clustered using an improved MajorClust algorithm to get a better cluster. This technique provides parameter-free clustering so that it automatically can produce an analysis report for the forensic investigators. The clustering results will be checked for anomalies based on a score that considers some factors such as the total members in a cluster, the frequency of the events in the log file, and the inter-arrival time of a specific activity. We also provide a graph-based visualization of logs to assist the investigators with easy analysis. Experimental results compiled on an open dataset of a Linux authentication log show that the proposed method achieved the accuracy of 83.14% in the authentication log dataset.

© 2017 Elsevier Ltd. All rights reserved.

## Introduction

Access controls limit what actions user can perform in a specific environment such as an application or operating system. In this research, we mainly consider the access controls in an operating system. For example, this log is available in `auth.log` under the `/var/log/` directory in a Debian-based Linux environment or `/var/log/secure` in RedHat family distributions. It provides a log of Pluggable Authentication Modules (PAM) that record the user's granted access (Geisshirt, 2007) and also secure shell (SSH) accesses, both failed and successful ones. This artifact is important evidence for the forensic investigators or security analysts to trace the attackers and analyze the incidents in a server.

For `auth.log` analysis, there is OSSEC (Open Source Host-based Intrusion Detection System Security), which also considers other log files to detect suspicious activities in a host (OSSEC). Some examples of OSSEC implementation to detect access violations and

multiple failed logins were presented in Cid (2009) while its complete guide is described in Hay et al. (2008). Other research on `auth.log` investigation was proposed by Sato and Yamauchi (2012). They first provided an architecture for securing log files and then proposed tampering detection of system logs including `auth.log`. In addition, Basin et al. gave a brief explanation of possible attacks and a simple examination using existing tools to check the traces for attempted attacks. As a result, we can see a trail in the security log such as `auth.log` (Basin et al., 2011).

One of the applications that record activities in `auth.log` is the SSH server. Detection of brute force attacks is critical because recent reports show this type of attack is frequently conducted where SSH service is widely targeted (Abdou et al., 2015). Instead of doing forensic analysis, most research have conducted a proactive mechanism to prevent the dictionary attack (Thames et al., 2008; Su, 2011). The authors presented a detection and defense architecture for SSH threats based on the authentication log. Another attempt to create a model for brute force profiling was proposed by Javed and Paxson where the model accounted the detection and user authentication failures (Javed and Paxson, 2013).

The application of text clustering in digital evidence has also been proposed to increase the relevance of the search results

\* Corresponding author. School of Engineering and Information Technology, Murdoch University, Australia.

E-mail addresses: [hudan@its.ac.id](mailto:hudan@its.ac.id), [studiawan@gmail.com](mailto:studiawan@gmail.com) (H. Studiawan), [c.payne@murdoch.edu.au](mailto:c.payne@murdoch.edu.au) (C. Payne), [f.sohel@murdoch.edu.au](mailto:f.sohel@murdoch.edu.au) (F. Sohel).

(Beebe and Clark, 2007). Several papers have tried to cluster the event log and detect its outliers, but the threshold for the anomaly scores was not sophisticated since the user must tune it manually to get the best results (Vaarandi, 2003, 2008).

In this paper, we consider SSH records as well as PAM authentication in `auth.log`. The proposed method identifies any log activities that are suspicious and labels them as outliers after the clustering phase. We propose a parameter-free algorithm to cluster this log and conduct anomaly detection, which refers to establish a baseline norm and detecting deviations from it. These deviations can be an attack or other strange events which need attention from the forensic analysts. However, the experiments focus on brute force password attacks as this is the most common type of access control violation (Abdou et al., 2015). The term parameter-free means that the clustering algorithm will run without any initial guess or any parameter supplied.

Most existing techniques need a manual input of the number of clusters. The classical k-Means and its variants still need the number of clusters as a parameter to run (Celebi et al., 2013) while hierarchical clustering requires a specific level at which to cut the generated tree or dendrogram (Murtagh and Contreras, 2012). Density-based clustering and its extensions are other alternatives, but they require the neighborhood size to be passed as a user-defined variable (Kriegel et al., 2011).

Another technique, namely MajorClust, is a parameterless method that clusters nodes based on its natural structure in the graph (Stein and Niggemann, 1999). MajorClust has been extended to a fuzzy version (Levner et al., 2007) and probability-based approach, so it can work well with several structures on the local scale (Niggemann et al., 2010). However, we will improve MajorClust by adding a condition when the processed node's cluster is not the same as the heaviest neighbor node's cluster, we will force it to follow the cluster of the heaviest one. Therefore, an event will stick to the most similar one instead of following the others, which are not so similar but has more edge weight aggregation. This will improve the performance of clustering so that it will be suitable with the log record. The next section will describe the proposed method based on refined MajorClust.

## Proposed method

We define some formal terms in this paper. An *event log* or a *log file* is a file that records activities from an application or services such as authentication log from SSH server application. A *record* is defined as a single entry in the event log. Then, an *event* is defined as a message in a record.

The proposed method is illustrated in Fig. 1. First, the raw log will be preprocessed and converted into a graph model. Second, we enter the clustering phase which consists of four steps:

1. Cluster graph using improved MajorClust (Phase 1);
2. Create new representation of graph based on multiple longest common substring (MLCS) to solve the overfitting problem (Phase 2);
3. Refine the new form of cluster by running improved MajorClust once again (Phase 3);
4. Produce clustering result using initial graph representation (Phase 4).

After that, we conduct anomaly calculation on the clustering results based on a score considering several properties that characterize each cluster. An estimated threshold for anomalousness is provided to detect this behavior automatically. The last step is to provide a visualization of the detected outliers. The details of each phase are described in the following subsections.

### Log preprocessing and proposed graph model

In this paper, the digital evidence to be investigated is the authentication log file (`auth.log`) in a Linux environment. It is assumed that the log file has not been tampered with or modified by the attackers. The full set of features for analysis are datetime, hostname, service or process name, process identifier (PID), and the event. The feature used in log preprocessing step is only the event since we focus on the content of the log and we do not consider other fields. It should be noted that we do not delete any fields from the log, rather only temporarily omit these features in the clustering process. For example, we need datetime property for calculating the anomaly score later. Other ignored properties include the hostname of the server, type of log and its process ID (PID), e.g., `sshd[2790]` and `CRON[2839]`. We filter all events, so it will produce only the unique one and attach the row log identifier to every unique event.

Furthermore, all numeric characters and stopwords in unique events are removed. We include some additional words that are not included in standard English stopwords but exist in the log record, e.g., *preauth*, *from*, *for*, *port*, *sshd*, *ssh*, and *root*. To identify non-standard stopwords, we conducted the experiments and checked the results of the graph clustering. If the results are not formed well, we analyze the string message and add the stopwords in order to increase the clustering accuracy. The added stopwords commonly appear in the string and if not removed it will send the different

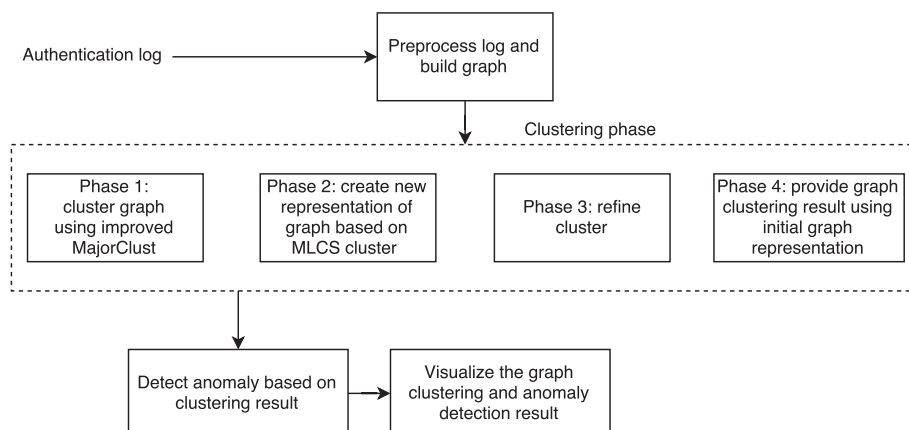


Fig. 1. Block diagram of proposed method.

messages into one cluster.

Moreover, a *tf* – *idf* (term frequency–inverse document frequency) procedure is implemented in every filtered line to produce its numerical representation. Note that in our case, document *d* refers to a single line of log record *l*. First, term frequency, *tf*, is not only the number of occurrences of term *t* in log *l* but it is normalized with the total number of terms in *l*, which is denoted as *len(l)* so that:

$$tf_{t,l} = \frac{tf}{len(l)} \quad (1)$$

Second, we need to calculate *inverse document frequency* of term *t*, *idf<sub>t</sub>*, as shown in the equation below:

$$idf_t = 1 + \log \frac{N}{df_t} \quad (2)$$

where *N* is the number of lines in the event log, and *df<sub>t</sub>* is the total number of lines in which the term occurs. We modify a basic formula of *tf<sub>t,l</sub>* and *idf<sub>t</sub>* from Manning et al. (2009) to conform with the real event log and avoid a zero value. Finally, we calculate *tf* – *idf* of term *t* as follows:

$$tf - idf_{t,l} = tf_{t,l} \times idf_t \quad (3)$$

The next step is to build a graph  $G = (V, E, w)$  where *V*, a set of vertices, represents each unique event, *E*, a set of edges, depicts relationship between two vertices where its weight, *w*, is the cosine distance to measure the similarity between two events in a single log record *l* as seen below (Manning et al., 2009):

$$w(l_1, l_2) = \frac{\vec{V}(l_1) \cdot \vec{V}(l_2)}{|\vec{V}(l_1)| |\vec{V}(l_2)|} \quad (4)$$

where the numerator denotes the dot product of the vector  $\vec{V}(l_1)$  and  $\vec{V}(l_2)$  which is represented as  $\sum_{i=0}^p \vec{V}_i(l_1) \vec{V}_i(l_2)$ . On the other hand, the denominator is the product of their Euclidean lengths and defined as  $\sqrt{\sum_{i=0}^d \vec{V}_i^2(l)}$  where *d* is the total number of the term calculated in a record *l*. This cosine similarity distance becomes the edge weight of two vertices. An example of the proposed graph model is shown in Fig. 2. We can see that every vertex is connected to other vertices except its zero distance since the edge is only created when the cosine similarity is greater than zero.

To provide a more clear illustration of cosine similarity in the event log, we give a step by step example of how this measurement is calculated between two events as follow.

**Step 1:** Two events from an event log

Dec 1 23:05:20 ip-172-31-27-153 sshd[28547]: Invalid user test from 192.208.179.82

Dec 1 23:36:42 ip-172-31-27-153 sshd[28578]: Invalid user admin from 187.76.79.142

**Step 2:** Preprocessing

*l<sub>1</sub>*: invalid user test

*l<sub>2</sub>*: invalid user admin

**Step 3:** Calculating term-frequency (*tf*)

Event	Invalid	User	Test	Admin
<i>l<sub>1</sub></i>	1	1	1	0
<i>l<sub>2</sub></i>	1	1	0	1

**Step 4:** Calculating normalized term-frequency (*tf*) based on Equation (1)

Event	Invalid	User	Test	Admin
<i>l<sub>1</sub></i>	0.3	0.3	0.3	0
<i>l<sub>2</sub></i>	0.3	0.3	0	0.3

**Step 5:** Calculating inverse document frequency (*idf*) based on Equation (2)

$$idf_{invalid} = 1 + \log 2/2 = 1$$

$$idf_{user} = 1 + \log 2/2 = 1$$

$$idf_{test} = 1 + \log 2/1 = 0.693$$

$$idf_{admin} = 1 + \log 2/1 = 0.693$$

**Step 6:** Calculating *tf* – *idf* based on Equation (3)

Event	Invalid	User	Test	Admin
<i>l<sub>1</sub></i>	0.3*1	0.3*1	0.3*0.693	0
<i>l<sub>2</sub></i>	0.3*1	0.3*1	0	0.3*0.693
<i>l<sub>1</sub></i>	0.3	0.3	0.208	0
<i>l<sub>2</sub></i>	0.3	0.3	0	0.208

**Step 7:** Calculating the numerator of cosine similarity

$$\begin{aligned} |\vec{V}(l_1)| \cdot |\vec{V}(l_2)| &= 0.3*0.3 + 0.3*0.3 + 0.208*0 + 0*0.208 \\ &= 0.09 + 0.09 + 0 + 0 \\ &= 0.18 \end{aligned}$$

**Step 8:** Calculating the denominator of cosine similarity

$$\begin{aligned} |\vec{V}(l_1)| &= \sqrt{0.3^2 + 0.3^2 + 0.208^2 + 0^2} \\ &= \sqrt{0.09 + 0.09 + 0.043 + 0} \\ &= \sqrt{0.223} \\ &= 0.472 \\ |\vec{V}(l_2)| &= \sqrt{0.3^2 + 0.3^2 + 0^2 + 0.208^2} \\ &= \sqrt{0.09 + 0.09 + 0 + 0.043} \\ &= \sqrt{0.223} \\ &= 0.472 \end{aligned}$$

**Step 9:** Calculating the cosine similarity based on Step 7, 8, and Equation (4)

$$\begin{aligned} w(l_1, l_2) &= \text{numerator/denominator} \\ &= 0.18/(0.472*0.472) \\ &= 0.18/0.228 \\ &= 0.789 \end{aligned}$$

*Event log clustering based on improved MajorClust*

The first phase to detect the anomaly is to cluster the generated graph using the MajorClust algorithm. For the sake of completeness, we include a brief technical description of MajorClust. For details of the algorithm, the reader is referred to Stein and Niggemann (1999). Initially, each node is attached to its own cluster. For each vertex, the procedure will accumulate the edge weight of the neighboring nodes as denoted below:

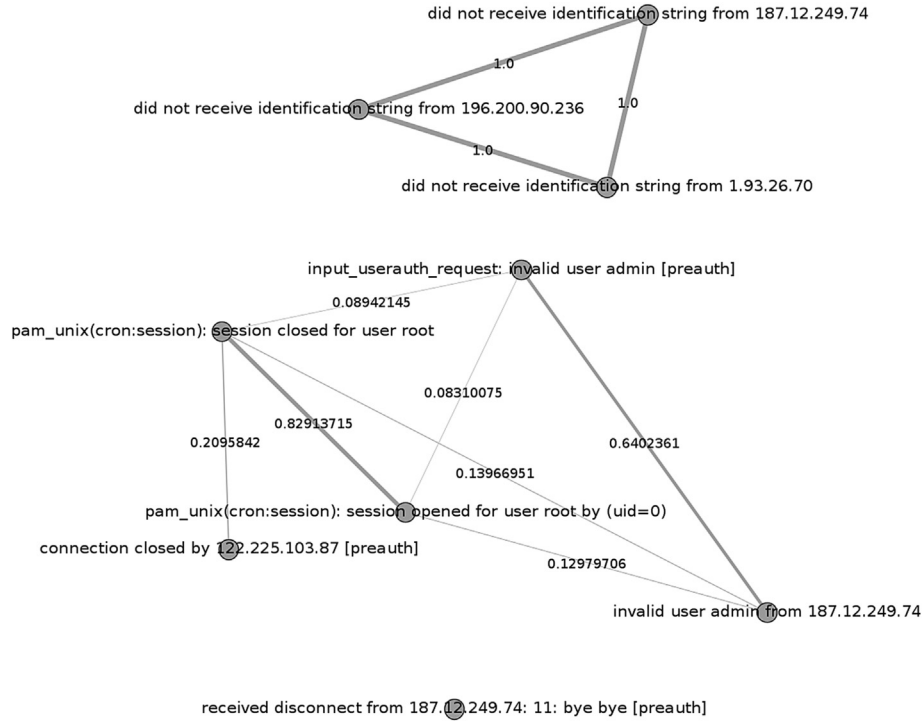


Fig. 2. Proposed graph model for authentication log.

$$c_i = \sum_{j=1}^n w(e_{ij}), \quad 0 \leq i \leq m \quad (5)$$

where  $c_i$  is the  $i$ th neighboring clusters,  $w(e)$  is the edge weight for each neighboring node in a particular cluster,  $m$  is the total of the neighboring cluster  $i$ , and  $n$  is the total member of vertices in the cluster  $c_i$ . The next step is assigning each vertex to a cluster, which has maximum weight aggregation and is defined as:

$$c^* = \operatorname{argmax}_i(c_i) \quad (6)$$

This process is continued until all possible combination checking for current and assigned clusters for each node is accomplished.

However, we identified a major drawback in the original MajorClust technique when we applied it in a `auth.log` file. It will fail to cluster a large number of nodes since it only depends on the accumulated weight of the neighboring vertices. A node will only depend on the aggregate weight of the neighbors although it is less relevant and does not consider the more similar cosine distance.

This shortcoming is illustrated in Fig. 3 and indicated by a large green node with label `pam_unix(cron:session): session closed for user root`. This node should create its own cluster since it is closer with another vertex with a more similar event.

Therefore, we improve the MajorClust by supplying the additional requirement that when the current cluster is not the same as the heaviest neighbor node, we force the operating node to follow the cluster of the heaviest one. The heavier the edge weight, the closer the distance between two events in the log record since it represents the cosine similarity. By deploying this improvement, an event will stick to the most similar one instead of following the others, which are not so similar but has more edge weight aggregation. We then refine Equation (6) as follows:

$$\mathbf{c} = \begin{cases} \operatorname{argmax}_i(c_i) & \text{if } c^* = c_i(h) \\ c_i(h) & \text{if } c^* \neq c_i(h) \end{cases} \quad (7)$$

where  $\mathbf{c}$  is the current cluster considering additional checking,  $h$  is a

neighbor node with the heaviest edge weight, and  $c(h)$  is the cluster label of  $h$ . The illustration of this improvement is depicted in Fig. 4. We can see that the node in the previous figure has created a new cluster as expected (indicated by vertices in red (in the web version)).

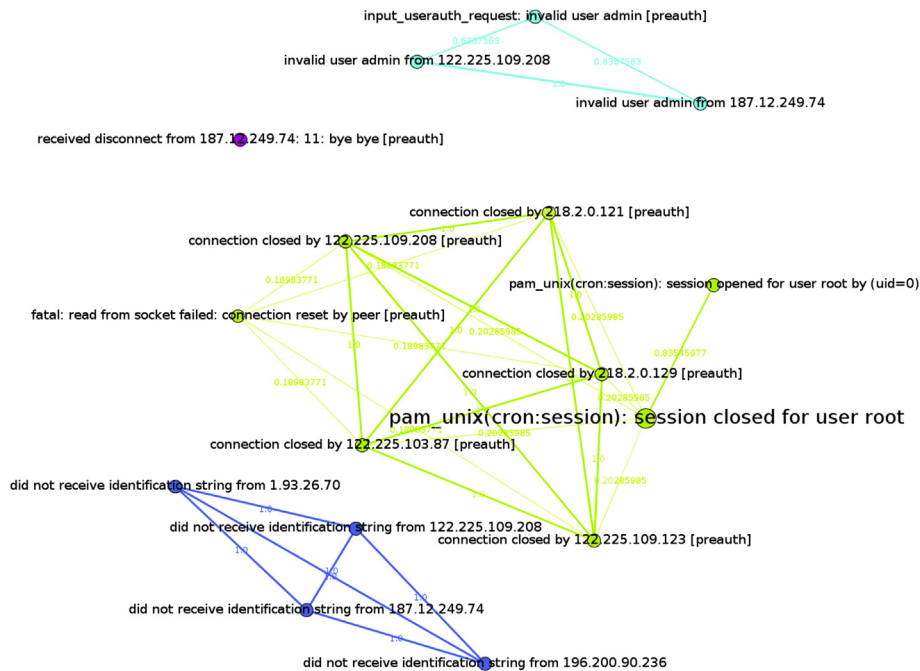
Nevertheless, this improvement triggers another problem. The generated cluster is overfitting since it produces many small clusters with only two or three vertices as shown in Fig. 5. This issue was found when the proposed procedure processed the first 500 records of the first-day log in the dataset. We solve this problem by running one additional round of MajorClust called the *refine cluster* phase and is described as follows. First, we represent each cluster as a single node. We maintain the start and end time of the overall event log for each cluster to be considered for the anomaly score calculation. The frequency of an event occurring in all nodes in the cluster and updated  $tf - idf$  is also counted.

To represent an event of a newly formed node, we need a string which reflects all events in a cluster. We can consider this as the longest common subsequence problem with many strings, or well-known as the multiple longest common substring technique. Solutions of this classical problem have been proposed in the recent literature (Yang et al., 2013, 2014). For each cluster  $\mathbf{c}$ , we will check for the most common substring occurring in the raw log. The terms “subsequence” and “substring” will be used interchangeably since the substring has the same meaning as a subsequence in our case, and we then describe MLCS concisely.

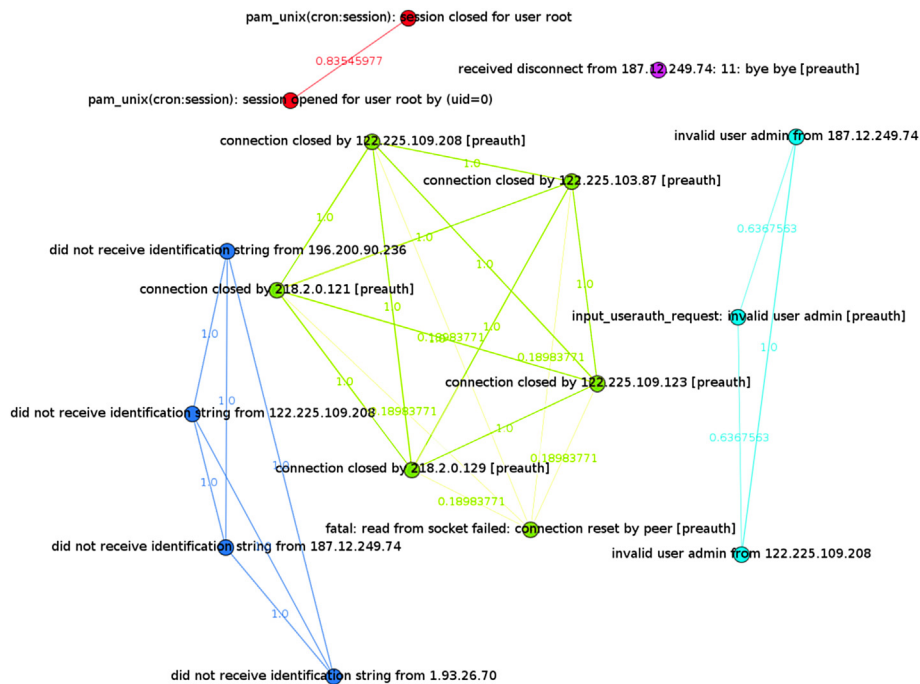
Let  $l = l_1 l_2 \dots l_p$  be a single log string in a particular cluster where the subscript number represents the index of a character in  $l$ ,  $s = s_{w_1} s_{w_2} \dots s_{w_q}$ ,  $p$  and  $q$  are their length respectively,  $s$  is called subsequence of  $l$ ,  $\operatorname{sub}(l, s)$ , if it meets the following:

$$\operatorname{sub}(l, s) = \begin{cases} 1 \leq v \leq q & \text{if } 1 \leq w_v \leq p \\ 1 \leq t < u \leq q & \text{if } w_t < w_u \end{cases} \quad (8)$$

Furthermore, let  $L = \{l_1, l_2, \dots, l_o\}$  be a set of raw events in the specific  $\mathbf{c}$ ,  $o$  is total number of logs within each  $\mathbf{c}$ , multiple longest common substring for set  $L$  is a sequence of  $s$  if and only if (i)  $s$  is the



**Fig. 3.** The drawback of MajorClust with more vertices.



**Fig. 4.** The improvement of MajorClust algorithm.

subsequence of  $l_i$  for  $1 \leq i \leq o$  and (ii)  $s$  is the longest one satisfying (i). We then apply MLCS to all clusters to find the longest substring of  $L$ .

After that, we run the MajorClust algorithm again using current graph composition. The result of the *refine cluster* phase is given in Fig. 6. After this phase is complete, we convert the generated graph back to its original nodes representation since we still maintain the initial vertices of the formed new node. This technique will provide us with a better clustering result as shown in Fig. 7. The output of

this phase is a complete graph where each node has its own cluster label property.

It should be noted that the algorithm converges smoothly. The reasons are two-fold. First, the clustering algorithm works on a node by node basis. When a node is processed, the node is flagged, and the algorithm moves to the next node and so on. Second, the proposed improvement of MajorCluster will ensure a node to have the same cluster as its heaviest edge neighbor. Using these two conditions, the algorithm will smoothly converge without any oscillation.





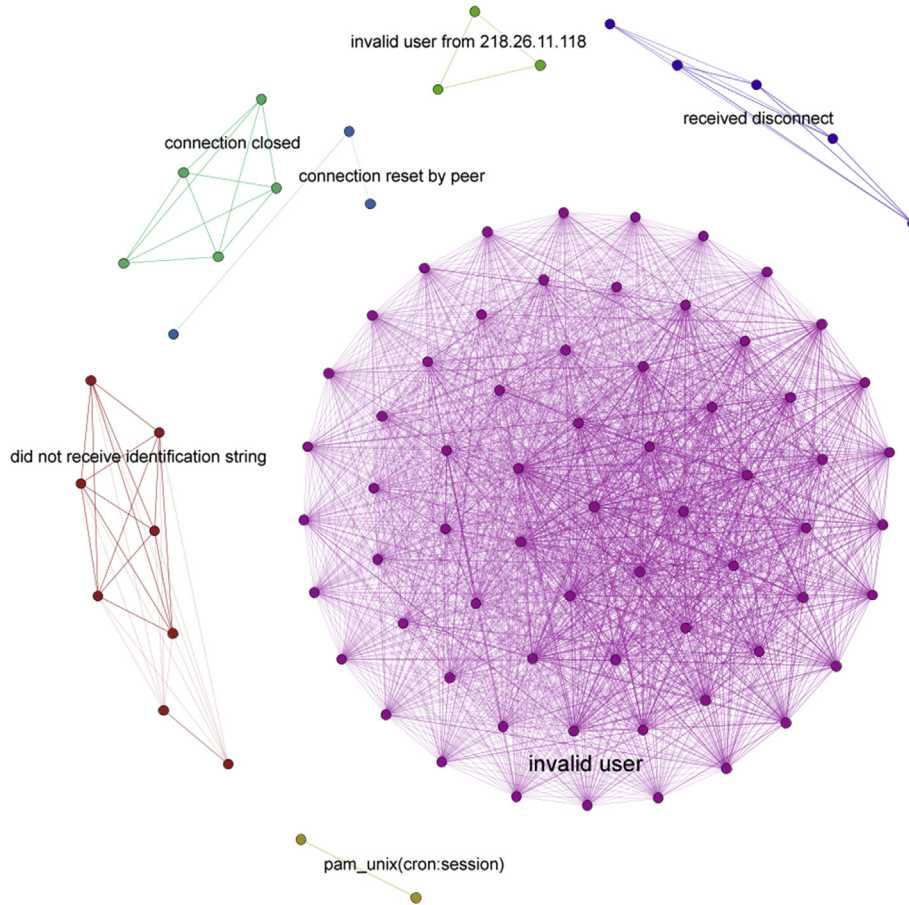


Fig. 7. The final clustering of improved MajorCluster algorithm.

*refine cluster* phase, respectively. In addition,  $n$  is the total member of vertices in the cluster  $c_i$ . The anomaly will occur when there are so many events occurring in the adjacent time while the normal event shows the opposite behavior. Formally, the inter-arrival rate of all events in a cluster,  $\mathbb{I}$ , is formulated as follows:

$$\mathbb{I} = \frac{\sum_{i=1}^n f_i}{t_n - t_1} \quad (10)$$

where  $t_1$  is the first time an event occurred,  $t_n$  is the last one, and both of them are measured in seconds. When the denominator produces a zero value, it means that there are several events in a second, and we set  $t_n - t_1$  with very small value to show that its arrival rate is very high in the short time frame.

Furthermore, the anomaly score per cluster  $a$  is calculated based on the formula below:

$$a = \frac{\mathbb{M} \times \mu_f}{\sum_{i=1}^z \mathbb{N}_i \times \sum_{i=1}^z \mathbb{F}_i} \times \mathbb{I} \quad (11)$$

where  $\mathbb{M}$  is total members in  $c$ ,  $\mathbb{N}$  is total nodes in a cluster,  $\mathbb{F}$  is total frequencies in  $c$ , and  $z$  is total clusters in the analyzed log data. We refer to Equations (9) and (10) for  $\mu_f$  and  $\mathbb{I}$ , respectively.

#### Estimation of anomaly threshold

To decide whether or not a cluster is an anomaly, we estimate a threshold to provide a recommendation for the forensic investigator or security analyst. We introduce an estimation for the anomaly parameter. There are three types of activities recorded in

the log file: low intensity, normal, and high intensity (Javed and Paxson, 2013; Zhong et al., 2011). Low intensity refers to a sporadic attack while high intensity usually related to a brute-force attacks. In addition, the normal activities will produce a non-suspicious log.

First, the anomaly score that is calculated before is normalized to  $[0,1]$  as below. This will make a number of scores to fall in the same range.

$$a'_i = \frac{a_i - \min(A)}{\max(A) - \min(A)} \quad (12)$$

where  $a_i$  is an anomaly score in a cluster  $i$ ,  $A$  is a set of anomaly scores for all clusters, and  $a'_i$  is a normalized score of  $a_i$ .

Next, we fit this normalized score to a quadratic equation since this equation fits the characteristics of the dataset where low and high intensity event will be defined as anomalies. The quadratic equation fits  $a'$  as axis to  $[0,1]$  and  $a''$  as ordinate to  $[0,1]$ . The higher  $a''$  score, the more normal the events.

$$a''_i = a_i'^2 + 4a'_i - 4 \quad (13)$$

Subsequently, we have a second normalization step for anomaly score  $a''$  to range  $[-1,1]$  where 0 is the threshold for the anomaly decision. The final anomaly score per cluster  $\alpha_i$  is depicted below where  $p = -1$  and  $q = 1$  as the range boundaries.

$$\alpha_i = p + \frac{(a''_i - \min(A'))(q - p)}{\max(A') - \min(A')} \quad (14)$$

If the  $\alpha_i$  score is less than 0, then a cluster  $i$  is set as an anomaly and vice versa. Thus, the user is not required to enter the parameter to calculate the anomaly score.

#### Visualization of access control anomaly

In order to help the investigators in analyzing and getting a better understanding of the security log, there are some methods available for log visualization. For example, Takada and Koike created Tudumi, a visualization tool for auditing `syslog`, `wtm`, and `su` based on layered concentric disks (Takada and Koike, 2002). For each layer, there were some notations such as spheres and cubes to represent the user and his activities. The treemaps model to display clustering result from the Simple Log file Clustering Tool (SLCT) (Vaarandi, 2003) of the event log was introduced and named LogView (Makanju et al., 2008). Another approach was using parallel coordinates to plot several logs, e.g., network, database, and `syslog`.

Elvis (Extensible Log Visualization) is a recent implementation to display an Apache log, `syslog`, and `auth.log` (Humphries et al., 2013). This technique was based on a custom organization (categorical and geographical) and log augmentation. Visual Filter enabled the security analyst to inspect the whole logs, creates a filter visually, performs navigation, and sub-selects the part of the log interactively (Stange et al., 2014). In addition, Trethowen et al. introduced VisRAID to visualize remote access logs, especially for intrusion detection purposes and focused on timeline visualization (Trethowen et al., 2015).

However, existing methods do not apply graph visualization to the assistance tool for log analysis. Naturally, since the clustering is based on a graph, the proposed visualization also heavily relies on this model. Therefore, we propose analysis and anomaly display of the authentication log based on graph visualization. We categorize the visualization into two types: static and dynamic. The static model displays the graph after the analysis is complete while the dynamic one provides live visualization when the analysis is running.

We use Gephi for the graph visualization for both types (Bastian et al., 2009). We then exploit the streaming server plugin (Panisson, 2013) and its Python client implementation (Panisson, 2014) to support a dynamic model. In the static mode, the visualization is displayed after the clustering and anomaly detection process is finished. The output of the forensic analysis is a graph file, i.e., dot file and then it is exported to Gephi. To get the intuitive display, we use two steps of the graph layout algorithm: first Force Atlas 2 (Jacomy et al., 2014) and then Fruchterman and Reingold (1991). These two algorithms have been natively integrated into Gephi.

Force Atlas 2 produces clustered nodes, but there are still many overlaps between them, Fruchterman Reingold will remove these obstacles, and provide a clearer layout. While in the dynamic mode, the visualization follows the analysis process starting from creating nodes, edges, graph clustering, refining the cluster, and anomaly detection. The result is clean and concise since the procedure has removed unnecessary edges that connect one cluster to the other one. For example, Figs. 2–7 are produced in the static mode.

## Experimental results and discussions

#### Functionality testing for SecRepo dataset

The dataset for the experiments is taken from the public and open Security Repository (SecRepo). It includes 86,839 lines of `auth.log` taken from November 30 to December 31, 2014, and mainly contains failed SSH login attempts (Sconzo, 2016). This kind of attack will be flagged as an anomaly in the whole log file. We utilize NetworkX to create and manipulate the graph (Hagberg

et al., 2008), Gephi as graph visualization tools (Bastian et al., 2009), a Gephi streaming server plugin (Panisson, 2013) and Gephi streaming client (Panisson, 2014) to dynamically draw a graph to Gephi, and Natural Language Toolkit (NLTK) to provide English corpus for log text preprocessing, especially in providing stopwords (Bird et al., 2009). We utilize the Python programming language version 2.7 to assemble all of these libraries and provide this log forensic tool.

Fig. 8 shows the initial clustering result of the first day in the dataset (November 30, 2014) while Fig. 9 illustrates the *refine cluster* phase. There are too many small clusters in the initial step and it is fixed in the next phase. The anomaly detection result is depicted in Fig. 10 where the red vertices represent the anomaly and the green represents the normal events. This figure is completed with the label of an event in the access control log. The outlier is measured when a cluster is on the *refine cluster* phase based on the estimated threshold.

We can see in Fig. 10 that the detected anomaly is an `invalid user` event as shown in the largest cluster. There are small clusters with similar events such as `invalid user tomcat`, `invalid user cms`, `invalid user dev`, and `invalid user google`. These clusters create their own cluster since the nodes inside them have the neighbor with the heaviest weight and do not follow the bigger cluster, which has less cosine similarity.

In addition, the `pam_unix(cron:session)` session is a daemon event run in the background so it is successfully categorized as normal. The other two nodes in green, i.e., `connection close by` and `connection reset by peer` is a false positive since these events will only happen when there are failed authorizations. The calculation of the accuracy and performance of the proposed method will be explained later in this section.

Fig. 11 is the full version of Fig. 10 which is converted back to its initial graph representation where it has 688 rows of the event log and is modeled into 173 nodes and 9694 edges. To produce a good display of anomaly detection, these visualizations are generated after two runs of the layout algorithm, i.e., Force Atlas 2, then followed by Fruchterman Reingold, and successfully implemented in the Gephi graph editor for both static and dynamic modes.

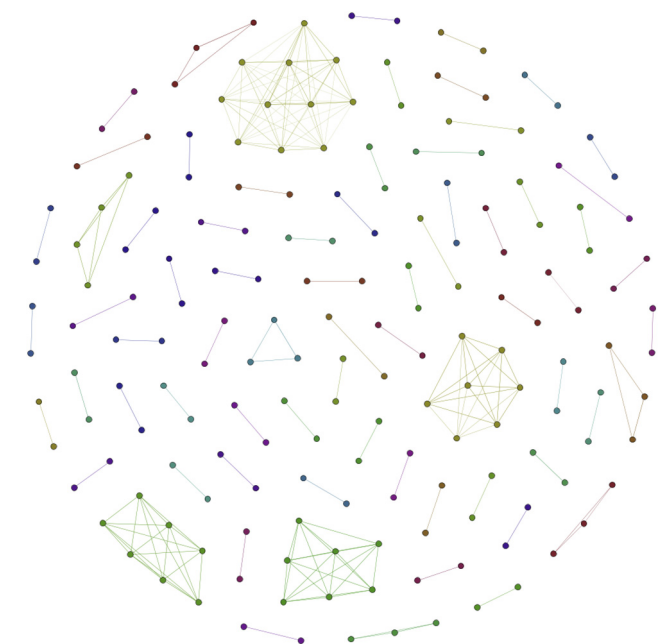


Fig. 8. Initial clustering result of the first day in the dataset (November 30, 2014).



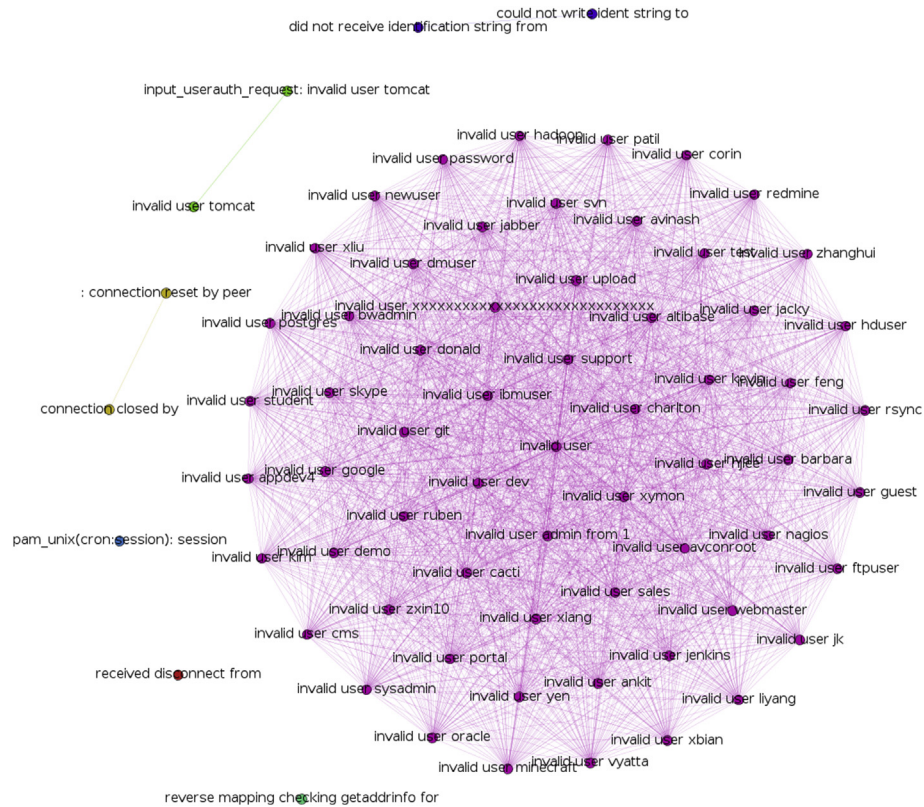


Fig. 9. The result of *refine cluster* phase from Fig. 8.

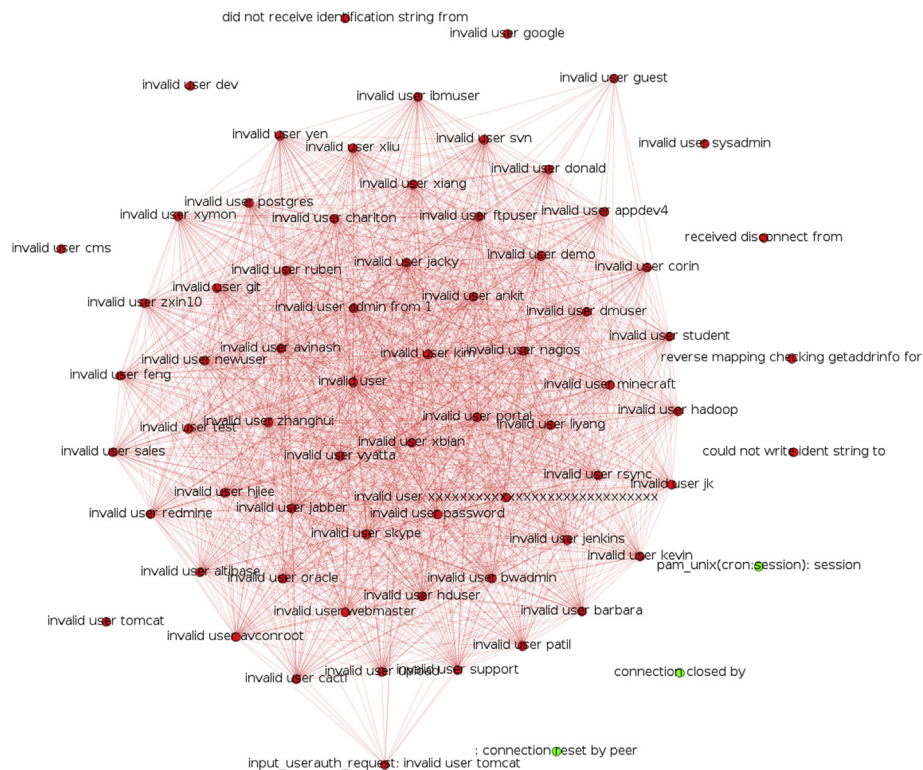


Fig. 10. Anomaly detection in *refine cluster* phase.

### Evaluation metrics

Furthermore, we use some standard measurements to evaluate the performance of the proposed method such as true positive ( $tp$ ), true negative ( $tn$ ), false positive ( $fp$ ), false negative ( $fn$ ), sensitivity, specificity, and accuracy. The sensitivity and specificity are defined in Equations (15) and (16), respectively (Manning et al., 2009):

$$\text{Sensitivity} = \frac{tp}{tp + fn} \quad (15)$$

$$\text{Specificity} = \frac{tn}{fp + tn} \quad (16)$$

While sensitivity is focused on measuring the improved MajorClust's ability to detect a normal event, the specificity and accuracy will provide an overall measure since it also includes true negative ( $tn$ ) as denoted below:

$$\text{Accuracy} = \frac{tp + tn}{tp + fp + fn + tn} \quad (17)$$

In this case, true negative refers to detected *anomaly* record is also defined as *anomaly* in the dataset. Before calculating  $tp$ ,  $tn$ ,  $fp$ , and  $fn$ , we label the dataset with *normal* and *anomaly* then compare it to the proposed method's decision.

### Comparison with existing methods

To compare the robustness of the improved MajorClust and proposed anomaly score, we evaluate other methods to the same testing dataset. We run the Simple Log Clustering Tool (`sclt`) version 0.05 (Vaarandi, 2003) and `LogCluster` version 0.03 (Vaarandi and Pihelgas, 2015) as the clustering-based anomaly detection in event logs. In addition, we also compare the proposed method with two rule-based anomaly detections, i.e., `OSSEC` version 2.8.3 (OSSEC; Cid, 2009; Hay et al., 2008) and `swatch` version 3.2.3 (Atkins, 2015). These two applications can be turned into forensic analysis tools since they were initially designed to monitor real-time log file. We also compare the proposed method with standard MajorClust and improved MajorClust without refine cluster phase. They will be combined with the proposed outlier detection mechanism since they do not have one.

### Parameter tuning for `sclt` and `LogCluster`

To objectively compare those methods, we have trained other tools' parameter to get their best accuracy. Therefore, the dataset was divided into two parts. The first part is November 30 to December 14 as the training set and the second one is December 15 to 31 as the testing set. The training set is used for parameter tuning while the testing set is for performance comparison as presented in Subsection Comparison results.

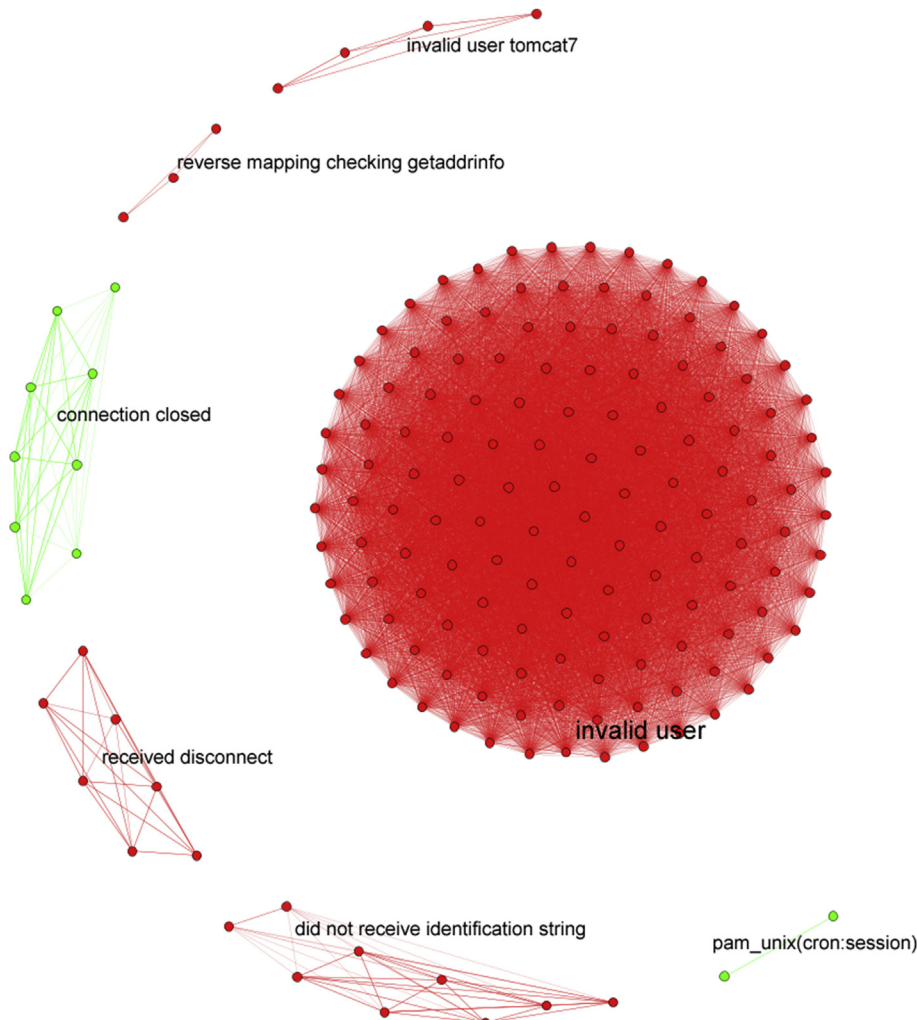


Fig. 11. Final result of the forensic analysis in the first day of authentication log (November 30, 2014).

The parameter `-support` or `-s` in `slct` and `-support` in `LogCluster` have been tuned from 10 to 100. Then we checked which value generates the best accuracy. Based on the experiments, parameter `-s` for `slct` is set to 100 which means that the threshold for the number of cluster member is 100 lines. As a consequence, log lines that do not belong to any cluster will be defined as an anomaly. The similar parameter is applicable to `LogCluster` and we set the option `-support` to 100 since it produces the best accuracy. The result of parameter tuning for `slct` and `LogCluster` is presented in Table 1.

The rules used in `OSSEC` and `swatch` are based on regular expressions. Since the `swatch` tool does not provide any standard rules, we set them by adding three most frequent malicious strings in the event log. They are ([iI]nvalid [uU]ser, Did not receive identification string, and reverse mapping). On the other hand, the `OSSEC` tool has provided the set of rules in an XML file. We followed these default rules.

### Comparison results

The comparison of the proposed technique and the state of the art is given in Table 2. As shown in Table 2, the best sensitivity value was achieved by `swatch`. This indicates that it has a very good capability to detect the normal activities in the log since we manually configure the regular expression rules that commonly occurred in the testing set. The proposed method achieved the best sensitivity and specificity rate of 70.59% and 82.21%, respectively. On the other hand, the other techniques produce lower specificity. The improved `MajorClust` is able to provide highest accuracy value (83.14%) while the others do not. The difference in accuracy is significant which indicates that existing methods are unable to properly detect the suspicious logs in the testing data.

We can see that the standard `MajorClust` without modifications does not work well and even it gives zero sensitivity. On the other hand, the `refine cluster` phase is hardly needed to improve the accuracy. Note that these two methods use the proposed anomaly score estimation.

The drawback of previous clustering-based anomaly detection is the assumption that the outlier is the data residing on a small

cluster. This idea does not work in the authentication log with a high number of attacks since these violations produce many records and can create the largest cluster in the analysis. Moreover, the rule-based anomaly detections contribute to the poor performance of anomalous discovery when the supplied regular expressions do not suit with the case. The default rule should be updated regularly by the network analysts in order to make sure all security breaches are detected. However, it can produce good true positive rates (sensitivity).

To supply the forensic investigator with more insight, we also enumerate the most frequent events when processing the `auth.log` file as presented in Table 3. The `invalid user` shows that there is an attempt from an unregistered user to login to the server. The event `received disconnect` means that an error occurred in the authentication process. The attacker is usually trying to conduct SSH brute force and the server decides to disconnect the connection.

There are warning messages such as `IP does not map back to the address` and `reverse mapping checking getaddrinfo failed` in the event. The system records these two events because the client's reverse DNS record does not match with the hostname used to identify the client's identity. Event `pam_unix(-cron:session)` means that the daemon activity is generated from the default configuration of Linux operating system. This record is written when the system checks the authentication mechanism and it happens every hour by default. From this analysis, we can infer that there are so many attempts to violate the access control, especially SSH, in the server.

### Experiment on the Kippo log

In this section, we present one more experiment to another dataset, i.e., Kippo log, to prove that the proposed method works well to another dataset with another attack. Kippo is an SSH honeypot and created in Python (Kippo). We install and configure Kippo on DigitalOcean droplet with public IP address so everyone including real attackers and botnet can reach this server. All the activities were recorded from 14 to 20 February 2017.

Beside the attack from the internet, we add one more activity that is usually done by the attacker before attempting penetration. The attack is a part of reconnaissance step to get banner or fingerprint of the SSH server version using the telnet application. This experiment shows that the proposed method performs well with 52.89% sensitivity, 92.77% specificity, and 87.35% accuracy.

**Table 1**

Parameter tuning `-s` for `slct` and `-support` for `LogCluster` (bold value indicates the best result).

Value of <code>-s</code> or <code>-support</code>	Accuracy for <code>slct</code> (%)	Accuracy for <code>LogCluster</code>
10	2.38	5.06
20	2.32	3.47
30	3.04	3.56
40	2.99	4.82
50	3.43	6.23
60	2.64	6.00
70	3.34	7.65
80	3.26	8.28
90	3.36	9.51
100	<b>3.62</b>	<b>10.23</b>

**Table 2**

Comparison of proposed technique and the other methods (bold value indicates the best result).

Methods	Sensitivity (%)	Specificity (%)	Accuracy (%)
Proposed method	70.59	<b>82.21</b>	<b>83.14</b>
Proposed method without <i>refine cluster</i> phase	88.24	50.25	52.43
Standard <code>MajorClust</code> (Stein and Niggemann, 1999)	0.00	78.08	73.61
<code>slct</code> (Vaarandi, 2003)	87.50	4.88	13.54
<code>LogCluster</code> (Vaarandi and Pihelgas, 2015)	73.21	17.94	23.73
<code>OSSEC</code>	25.11	30.12	29.59
<code>swatch</code> (Atkins, 2015)	<b>100.00</b>	54.13	58.94

**Table 3**

The most frequent events in the authentication log.

No	Event	Frequency
1	Invalid user	12,743
2	Received disconnect	11,464
3	<code>pam_unix(cron:session)</code>	2,708
4	IP does not map back to the address	1,328
5	Reverse mapping checking getaddrinfo failed	1,216



## Conclusion and future works

We propose an improved MajorClust algorithm to cluster the authentication log as the basis of anomaly detection of suspicious activities related to access control violations. The experimental results show that the proposed method is able to provide assistance to the forensic investigator, security analyst, or system administrator in inspecting and visualizing the log file. In the experiment, it achieves 70.59% of sensitivity, 82.21% of specificity, and 83.14% of accuracy.

In future, we would like to improve the similarity measure between two log records using semantic analysis since it considers the meaning, context, and linguistic aspect while  $tf-idf$  only calculates the frequency of terms. This measurement is expected to increase the accuracy of the anomaly detection phase. The clustering of the log file is still an open challenge for good performance of the analysis since the improved MajorClust, in some circumstances, failed to produce a precise cluster. In addition, the proposed method can be deployed to other types of event logs such as syslog, web server log, web proxy log, and many more.

Since the log data has been converted to graph data structure, we will add tamper detection using graph watermarks in the future. This method works by adding a subgraph to the main graph completed with a graph key and the instance's key (Zhao et al., 2015). The instance can be another process which saves the log files and it can reside in the same server as event log or in other servers. Watermarks guarantee the integrity of the log files and also provide a leak detection, so the suspected instance will be easily detected.

It is also worth noting that the runtime of the authentication log forensic analysis can still be improved. Some research in the digital investigation area has considered this problem (Quick and Choo, 2014). We plan to implement the parallel version of the proposed method to increase the processing time so that the authorities will receive the result and report it immediately. This is one of our priorities since we believe that the evidence not only originates from a single host but can also be obtained from a cluster environment.

## Acknowledgments

This work is supported by the Indonesia Lecturer Scholarship (BUDI) from Indonesia Endowment Fund for Education (LPDP). This scholarship is a collaboration between the Ministry of Research, Technology, and Higher Education with the Ministry of Finance, Republic of Indonesia. We also thank two anonymous reviewers who have provided the constructive comments to this manuscript.

## References

- Abdou, A., Barrera, D., Oorschot, P.C.V., 2015. What lies beneath? Analyzing automated SSH brute-force attacks. In: Proceedings of the 8th International Conference on Passwords, pp. 72–91.
- Aggarwal, C.C., 2013. Outlier Analysis. Springer.
- Atkins, T., 2015. Swatch: Simple log watcher. URL <https://sourceforge.net/projects/swatch/>.
- Basin, D., Schaller, P., Schlöpfer, M., 2011. Logging and log analysis. In: Applied Information Security. Springer-Verlag, Berlin Heidelberg, pp. 69–80.
- Bastian, M., Heymann, S., Jacomy, M., 2009. Gephi: an open source software for exploring and manipulating networks. In: Proceedings of International AAAI Conference on Web and Social Media, pp. 361–362.
- Beebe, N.L., Clark, J.G., 2007. Digital forensic text string searching: improving information retrieval effectiveness by thematically clustering search results. Digit. Investig. 4, 49–54.
- Bird, S., Klein, E., Loper, E., 2009. Natural Language Processing with Python. O'Reilly Media, Inc.
- Celebi, M.E., Kingravi, H.A., Vela, P.A., 2013. A comparative study of efficient initialization methods for the k-means clustering algorithm. Expert Syst. Appl. 40 (1), 200–210.
- Cid, D.B., 2009. Log Analysis for Intrusion Detection, Tech. Rep.
- Fruchterman, T.M.J., Reingold, E.M., 1991. Graph drawing by force-directed placement. Softw. Pract. Exp. 21 (11), 1129–1164.
- Geissshirt, K., 2007. Pluggable Authentication Modules. Packt Publishing, Birmingham, UK.
- Hagberg, A., Schult, D., Swart, P., 2008. Exploring network structure, dynamics, and function using NetworkX. In: Proceedings of the 7th Python in Science Conference, pp. 11–15.
- Hay, A., Cid, D.B., Bray, R., 2008. OSSEC Host-based Intrusion Detection Guide. Syngress Publishing.
- Humphries, C., Prigent, N., Bidan, C., Majorczyk, F., 2013. ELVIS: extensible log visualization. In: Proceedings of the 10th Workshop on Visualization for Cyber Security, pp. 9–16.
- Jacomy, M., Venturini, T., Heymann, S., Bastian, M., 2014. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. PLoS ONE 9 (6), 1–12.
- Javed, M., Paxson, V., 2013. Detecting stealthy, distributed SSH brute-forcing. In: Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security, pp. 85–96.
- Kippo, S.S.H. HoneyPot. URL <https://github.com/desaster/kippo>.
- Kriegel, H.-P., Kröger, P., Sander, J., Zimek, A., 2011. Density-based clustering. Wiley Interdiscip. Rev. Data Min. Knowl. Discov. 1 (3), 231–240.
- Levner, E., Pinto, D., Rosso, P., Alcaide, D., Sharma, R.R.K., 2007. Fuzzifying clustering algorithms: the case study of majorclust. In: Proceedings of the 6th Mexican International Conference on Artificial Intelligence, pp. 821–830.
- Loureiro, A., Torgo, L., Soares, C., 2004. Outlier detection using clustering methods: a data cleaning application. In: Proceedings of KDDNet Symposium on Knowledge-based Systems for the Public Sector.
- Makanju, A., Brooks, S., Zincir-Heywood, A.N., Milios, E.E., 2008. LogView: visualizing event log clusters. In: Proceedings of the 6th Annual Conference on Privacy, Security and Trust, pp. 99–108.
- Manning, C.D., Raghavan, P., Schütze, H., 2009. An Introduction to Information Retrieval. Cambridge University Press.
- Murtagh, F., Contreras, P., 2012. Algorithms for hierarchical clustering: an overview. Wiley Interdiscip. Rev. Data Min. Knowl. Discov. 2 (1), 86–97.
- Niggemann, O., Lohweg, V., Tack, T., 2010. A probabilistic majorclust variant for the clustering of near-homogeneous graphs. In: Proceedings of the 33rd Annual German Conference on Artificial Intelligence, pp. 184–194.
- OSSEC, Open Source HIDS Security. URL <http://ossec.github.io/>.
- Panisson, A., 2013. Gephi streaming plugin. URL <https://marketplace.gephi.org/plugin/graph-streaming/>.
- Panisson, A., 2014. Python client for graph streaming on Gephi. URL [https://github.com/panisson/pygephi\\_graphstreaming](https://github.com/panisson/pygephi_graphstreaming).
- Quick, D., Choo, K.K.R., 2014. Impacts of increasing volume of digital forensic data: a survey and future research challenges. Digit. Investig. 11 (4), 273–294.
- Sato, M., Yamauchi, T., 2012. VMM-based log-tampering and loss detection scheme. J. Internet Technol. 13 (4), 655–666.
- Satoh, A., Nakamura, Y., Ikenaga, T., 2015. A flow-based detection method for stealthy dictionary attacks against secure shell. J. Inf. Secur. Appl. 21, 31–41.
- Sconzo, M., 2016. SecRepo.com: Security Data Samples Repository. URL <http://www.secrepo.com/>.
- Stange, J.-E., Dörk, M., Landstorfer, J., Wettach, R., 2014. In: Proceedings of the 11th Workshop on Visualization for Cyber Security, pp. 41–48.
- Stein, B., Niggemann, O., 1999. On the nature of structure and its identification. In: Proceedings of the 25th International Workshop on Graph-theoretic Concepts in Computer Science, pp. 122–134.
- Su, Y.-N., 2011. Developing the upgrade detection and defense system of SSH dictionary-attack for multi-platform environment. iBusiness 03 (01), 65–70.
- Takada, T., Koike, H., 2002. Tudumi: information visualization system for monitoring and auditing computer logs. In: Proceedings of the International Conference on Information Visualisation, pp. 570–576.
- Thames, J.L., Abler, R., Keeling, D., 2008. A distributed active response architecture for preventing SSH dictionary attacks. In: Proceedings of the 2008 IEEE SoutheastCon, pp. 84–89.
- Trethowen, L., Anslow, C., Marshall, S., Welch, I., 2015. VisRAID: visualizing remote access for intrusion detection. In: Proceedings of the 20th Australasian Conference, pp. 289–306.
- Vaarandi, R., 2003. A data clustering algorithm for mining patterns from event logs. In: Proceedings of the 2003 IEEE Workshop on IP Operations and Management, pp. 119–126.
- Vaarandi, R., 2008. Mining event logs with SLCT and LogHound. In: Proceedings of the IEEE Network Operations and Management Symposium, pp. 1071–1074.
- Vaarandi, R., Pihelgas, M., 2015. LogCluster - a data clustering and pattern mining algorithm for event logs. In: Proceedings of the 11th International Conference on Network and Service Management, pp. 1–7.
- Yang, J., Xu, Y., Sun, G., Shang, Y., 2013. A new progressive algorithm for a multiple longest common subsequences problem and its efficient parallelization. IEEE Trans. Parallel Distrib. Syst. 24 (5), 862–870.
- Yang, J., Xu, Y., Shang, Y., Chen, G., 2014. A space-bounded anytime algorithm for the multiple longest common subsequence problem. IEEE Trans. Knowl. Data Eng. 26 (11), 2599–2609.
- Zhao, X., Liu, Q., Zheng, H., Zhao, B.Y., 2015. Towards graph watermarks. In: Proceedings of the 2015 ACM Conference on Online Social Networks, pp. 101–112.
- Zhong, Y., Yamaki, H., Takakura, H., 2011. A grid-based clustering for low-overhead anomaly intrusion detection. In: Proceedings of the 2011 5th International Conference on Network and System Security, pp. 17–24.