

Holmes: An Efficient and Lightweight Semantic Based Anomalous Email Detector

Peilun Wu*, Fan Yan[†], Hui Guo[§]

Feng Qian**, Quanzhong Zhan^{||} and Chao Zhang[¶]

School of Computer Science and Engineering, University of New South Wales (UNSW), Sydney[§]*

Department of Cyber Security Capability, Sangfor Technologies Inc., Shenzhen[†]*

*Information Center, Ministry of Water Resources, Beijing^{¶||**}*

*Email: *wupeilun@sangfor.com.cn, [†]yanfan@sangfor.com.cn, [§]h.guo@unsw.edu.au*

***qianfeng@mwr.gov.cn, ^{||}zqz@mwr.gov.cn, [¶]zhangchao@mwr.gov.cn*

Abstract—Email threat is a serious issue for enterprise security. The threat can be in various malicious forms, such as phishing, fraud, blackmail and malvertisement. The traditional anti-spam gateway often maintains a greylist to filter out unexpected emails based on suspicious vocabularies present in the email’s subject and contents. However, this type of signature-based approach cannot effectively discover novel and unknown suspicious emails that utilize various up-to-date hot topics, such as COVID-19 and US election. To address the problem, in this paper, we present Holmes, an efficient and lightweight semantic based engine for anomalous email detection. Holmes can convert each email event log into a sentence through word embedding and then identify abnormalities based on those translated sentences. We found that, in an enterprise environment, there is a stable relation between sender and receiver, but suspicious emails are commonly from unusual sources and present a rare relation, and hence can be detected through the rareness selection. We evaluate the performance of Holmes in a real-world enterprise environment, where around 5000 emails are sent/received each day. In our experiments, Holmes shows a high capability to detect email threats, especially those that cannot be handled by the enterprise anti-spam gateway. It is also demonstrated through our experiment that Holmes outperforms other popular anomalous email detection tools.

Index Terms—spam detection, novelty detection, machine learning, intrusion detection, fraud detection, phishing, malvertisement.

1. Introduction

Though the instant messaging software, such as Facebook and WeChat, has gained increasing popularity, the email service is still indispensable for enterprises. Since the email service is a public-facing application, it can be targeted by the hacker as an easy entrance to the internal network. Based on our observations, fraud, malvertisement and spread-phishing are the main email threats frequently received by enterprise users. These emails use deceptive subjects to pretend and hide themselves. Usually, malware infected attachments or malicious URLs are embedded in

the email body to spoof recipients for further action. Once the attachment is downloaded or a link is clicked, the recipients’s system is compromised or the confidential information is leaked [1].

To alleviate the problem, an enterprise often deploys some anti-spam gateways to filter out unexpected emails. However, the associated techniques for spam detection, such as greylist and subject analysis, cannot effectively discover novel and unknown email threats that are elaborately constructed by utilizing various current hot topics, such as COVID-19, US election. These unknown threats can easily bypass the anti-spam gateway and successfully permeate the target system, leading to a series of damaging consequences, such as administrator account theft, database attack and financial blackmail.

In this paper, we introduce a novel artificial intelligence based anomalous email detector, Holmes, that can effectively tackle the challenges that we mentioned above. Holmes combines word embedding with novelty detection to discover anomalous behaviours from a high volume of mirrored SMTP traffic in a large-scale enterprise environment. To improve the result interpret-ability, we trace the real source IP addresses of suspicious emails in line with their geographical positions and further visualize the correlated relations in a directed-force graph. Our contributions are summarized as follows:

- We propose an efficient and lightweight semantic based anomalous email detector, Holmes, which can not only effectively discover new email threats but also maintain a low false positive rate in a real-world environment.
- Different from others detectors that usually require to examine the email bodies, Holmes can simply discover anomalies based on the email headers, which significantly reduces the cost of resource consumption and avoid the need to access email bodies (a sensitive security issue).
- We demonstrate the correlated relations of detected suspicious emails via graph visualization and show that the attacker portrait (based on their geographical positions) is in line with the cyber threat intelligence

provided.

- We evaluate Holmes with a commercial anti-spam gateway deployed in a real-world enterprise environment. Holmes not only can accurately detect those email threats that have been blocked by the anti-spam gateway, but also can discover a large number of email threats that have successfully escaped from the gateway.
- We also compare Holmes with several commercial email detectors offered by different security companies in VirusTotal [2], which shows that Holmes outperforms those detectors with a very high detection rate on a set of malicious emails.

The remainder of the paper is structured as follows. We begin with a brief discussion of some related work on email detection in Section 2. We then in Section 3 introduce the proposed semantic based anomalous email detector, Holmes. In Section 4, we present our evaluation results of Holmes and several commercial security products; the demonstration of how the visualization can be used to reconstruct the attack stories is also given in this section. The enhancements on Holmes for the real world implementation is given in Section 5. The related work and the future work are discussed in Section 6, and the paper is concluded in Section 7. As an add-on section, we append some extra discussions at the end of this paper.

2. Preliminary Knowledge

Anomalous emails can be classified into external threats and internal threats in accordance with MITRE ATT&CK Matrix [3]. External threats are the emails sent from external sources, whereas the internal threats are the emails sent from legitimate users within an organization, but whose email accounts have been stolen and used for the lateral movement attack. Most of previous research mainly focus on one specific threat type, such as URL-based lateral phishing [4], [5], [6], [7] or phishing web pages from search engine in a large-scale cyberspace [8]. There are still many open questions and unsolved challenges that need to be addressed holistically. Some issues and the existing solutions are presented below.

No Built-In Authentication in SMTP. The lack of a native authentication mechanism inside the SMTP service presents a security loophole to attackers. Attackers can easily forge the email header by pretending to be someone the recipient knows or from a business the recipient has a relationship with, so as to spoof recipients and avoid spam block lists [9]. To address the problem, several frameworks, such as SPF [10] (Sender Policy Framework), DKIM [11] (Domain Key Identified Mail) and DMARC [12] (Domain-Based Message Authentication, Reporting, and Conformance), have been developed to incorporate authentication into the email system. However, these designs are still not very effective in terms of implementation. When integrating authentication into the mail system

with the component-based software design, there are inconsistency issues between the software components offered by different parties [13], such as the incompatibility of mail forwarding servers, which allows numerous email threats escape detection.

Lack of Sensitivity to Unknown Variations. The unreliability of SMTP leaves email threats to have evolved with many variations, which are difficult to be discovered by the traditional security products. We have evaluated several malicious email detection modules within our internal security products, which use pattern matching of attack signatures for anomaly detection. None of them can discover the crafted phishing emails that utilize business-related content to pretend themselves look normal for evasion. We also have used the crafted phishing samples collected from our real-world hunting to evaluate the detection rate of security products offered by our competitors; Nevertheless, the evaluation result also shows their low sensitivity to unknown threats – in fact, all testing samples can successfully escape from the detection of 60+ engines in VirusTotal (Enterprise Service). This kind of low ability of detecting unknown attack variations has motivated the security community to turn to AI-based methods for anomaly detection.

High False Positive Rate. The research on anomaly detection for cyber threat hunting has been discussed for decades. The main concern on applying machine learning for anomaly detection is the significant false positive rate (FPR). Even though new designs are continuously proposed aiming for improvement [14], [15], [16], [17], [18], they were rarely evaluated in a real-world working environment, let alone put into use in commercial systems.

High Cost and Performance Bottleneck. The imbalance between the cost of data collection and the performance of algorithmic consumption is a significant challenge for most of the AI-based detectors. Though the complexity of AI computing algorithms has been constantly improved, most AI modules still require large computing and storage resources, which makes the existing attack detectors not easy to use and very slow to response attacks [19], [20], [21]. Furthermore, the detectors that use supervised machine learning require the labeled input data records and often need to be retrained once their performance begins to degrade, which also makes the machine learning ineffective for detection automation.

Lack of Provenance Analysis. So far few detectors have considered to integrate the provenance analysis within the detection mechanism. We believe provenance analysis is an important and enabling component in malicious email detection. Provenance analysis [22], [23], [24], [25] can reveal the attack story and the detail of attacker portrait behind the email, such as (1) where the email is from, (2) who the real sender is, (3) how the malicious shellcode executes, (4) what the potential correlations between malicious events

are. The above information is important for the security team to analyze the attack techniques, tactics and procedures (TTPs) and further assist the security experts to identify the individual attackers or organizations.

3. Holmes - Anomalous Email Detector

To address the above challenges, we introduce an efficient and lightweight semantic oriented anomalous email detector, Holmes, that can detect an email attack by analyzing the emails's recipient/subject which is available in the email header of SMTP.

Holmes consists of four main functions: 1) Word Embedding, 2) Novelty Detection, 3) Rareness Selection, and 4) Correlation Graph Analysis. For each email, Holmes takes its header and converts it into a numeric presentation through word embedding. The numeric data records are then input to the machine learning unit (Novelty Detection). The unit generates a list of novel emails, which are, in turn, processed by the rareness selection procedure to narrow down the detection targets. The detected results are finally presented in a human readable format and the correlations of the related email attacks are also pictured with a graph.

Holmes was originally written in Python with only 52 lines code. Based on the run-time analysis, Holmes can complete the entire detection in less than 73 seconds with 127 MB memory consumption on around 700 MB datasets (one day SMTP records) in a CentOS virtual server. We open source the Python code of the main detection functions in this paper in a hope to assist researchers to evaluate and reuse Holmes in their future research. The design of each function is described below.

```

1  def Doc2Vec(self, feature):
2      """
3      :param feature: SMTP features
4      :return: word vectors
5      """
6      documents = [TaggedDocument(doc, [i]) for
7 i, doc in enumerate(feature)]
8      model = Doc2Vec(vector_size=40, min_count
9 =2, epochs=40)
10     model.build_vocab(documents)
11     model.train(documents, total_examples=
12 model.corpus_count, epochs=model.epochs)
13     return model.docvecs.vectors_docs

```

Listing 1. Doc2Vec

3.1. Word Embedding

Since the email textual header information cannot be directly used for machine learning, **how to effectively represent textual data to the machine understandable** is important. Most algorithm engineers use OneHotEncoder [26], [27], [28] or OrdinalEncoder [29], or bag-of-words (BOW) [30], which can be simply implemented by the open-sourced library Scikit-Learn [31]. However, the three methods are not able to effectively memorize and maintain the data semantic correlations either in temporal or in spatial dimension.

Header Feature	Example
smtp.srcIp	185.156.172.29
direction	Inbound
smtp.dstIp	10.1.128.31
srcIp.country	United States
header.subject	Urgent! Secure your Account!
header.to	IT Center security@xxx[.]gov.xx
fileName	xxx.xls.doc.zip
User-agent	Microsoft Outlook Express 2.0

Figure 1. SMTP Header Feature

To address the problem, we use paragraph vector (Doc2Vec) [32] for the conversion, as detailed by the Python code in Listing 1. Compared to other conversion methods, Doc2Vec is able to better keep the semantics of the words or more formally the distances between the words, which can be of variable-length ranging from sentences to documents.

Doc2Vec is a semi-supervised learning algorithm. Its input is unlabeled but what will be learned is specified/supervised. In our code, the inputs are email headers and what to be learned are the features in the header, as shown in Fig. 1. Besides of some basic attributes such as subject, header.from or user-agent, which are often forged by hackers, we design two additional features that can also be used to help identify anomalies: the direction of email (direction) and the country of source IP address (srcIp.country).

The code in Listing 1 converts each email event to a feature vector. The feature vectors are then used for novelty detection.

3.2. Novelty Detection

Anomalous emails are usually unknown and novel. Their behaviors often deviate from the trace of normal activities. We use Local Outlier Factor (LOF) [33] to discover those emails with abnormal behaviors, as given in Listing 2 and Listing 3.

```

1  def local_outlier_factor(self, train_feature,
2      test_feature):
3      """
4      :param train_feature: training data
5      :param test_feature: testing data
6      :return: decision scores
7      """
8      LOF = LocalOutlierFactor(n_neighbors=20,
9 novelty=True, contamination=0.5)
10     LOF.fit(train_feature)
11     return list(LOF.decision_list(listing 2
12 on_function(test_feature)))

```

Listing 2. Local Outlier Factor (LOF)

The LOF algorithm shown in Listing 2 can learn the feature vectors of historical emails (i.e. the train_feature dataset in the code) then provide the outlier score for newly seen emails (from the test_feature dataset). There are some compelling advantages of applying LOF for novelty

detection: (1) It allows to train learning model on the data with contamination; (2) It has low computing complexity and can be used for online-learning, hence avoiding the performance degradation and the cost of retraining; (3) It is not sensitive to fine-tuning, which is good to the effectiveness and stability of parametric learning.

```

1  def novelty_analysis(self, factor, test_feature)
2  :
3  :param factor: decision scores of LOF
4  :param test_feature: testing data
5  :return novel/unseen samples
6  """
7  threshold = 0
8  outliers = []
9  novelty = []
10 for score in factor:
11     if score < threshold:
12         outliers.append(factor.index(score))
13
14 for index in outliers:
15     novelty.append(test_feature[index])
16 return novelty

```

Listing 3. Novelty Analysis

The decision scores from the LOF code can be negative and positive. The negative values indicate the abnormalities and the positive values indicate the normal behaviours. We regard any vector with a score smaller than a threshold is associated with an anomalous email, which is traced by the its index in the dataset, as shown in Listing 3.

3.3. Rareness Selection

If we consider the relation of sender and recipient in emails, anomalous emails are often associated with a weak relation. Emails sent from the hacker to the same recipient are often very rare. We therefore can further narrow down the malicious emails based on the sender-recipient relation – those abnormal emails with a weak sender-receiver relation will be selected as the final detection result, as described in Listing 4.

In our design, the relationship is measured based on the combination of a set of email features: source IP (src_ip), direction, sender (mail_from), and receiver (mail_to). For a strong sender-receiver relation, there should be many emails of the same IP-direction-sender-receiver value. Therefore, we count emails for different IP-direction-sender-receiver values and select those that have a low count value (smaller than a threshold) as malicious emails.

```

1  def rareness_selection(self, focus_data):
2  :param data: unseen samples
3  :return rare emails
4  """
5
6  relation = {}
7  rareness = []
8  for each_data in focus_data:
9      src_ip = each_data[0]
10     mail_from = each_data[4]
11     mail_to = each_data[6]
12     direction = each_data[1]

```

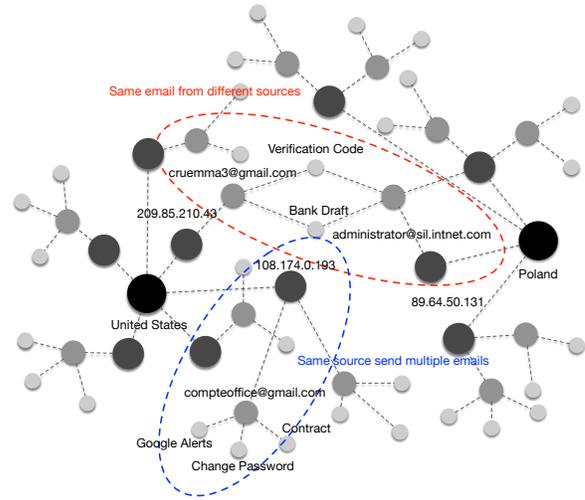


Figure 2. Correlation Graph Analysis

```

13     each_relation = src_ip + direction +
14     mail_from + mail_to
15     relation.setdefault(each_relation, [])
16     .append(each_data)
17     for k, v in relation.items():
18         if len(v) == 1:
19             rareness.extend(v)

```

Listing 4. Rareness Selection

With the word embedding, LOF novelty detection and rareness selection, Holmes can effectively discover anomalous emails.

3.4. Correlation Graph Analysis

Most of prior research overlooked a problem: **what is the relation within the anomalies?** Lack of an effective solution significantly increases the load of security analysts, blurs the attacker portraits, and further makes the provenance analysis hard. To address the problem, we introduce a correlation graph analysis (CGA) module to improve the clarity of attacker portrait descriptions by correlating different anomalous events.

CGA is a directed-force graph [34] and in our design, each node consists of the selected header features: country, srcIp, sender and subject. The directed graph enforces the nodes that have dense connections come closer but separates the nodes if they do not or have sparse connections. The graph depicts the similarity of different anomalies (such as the same srcIp, same subject or same sender) and centralizes the cluster in line with their geographical locations, hence significantly improving the interpret-ability of provenance analysis.

Fig 2 demonstrates the visualization result of CGA, where two clusters (one in red and one in blue) highlight the connected components that are centralized in accordance with the country of srcIp. The blue cluster shows that the same malicious email but sent from different sources,

and the red cluster reveals the same source sends multiple different malicious emails.

The CGA module can be used to generate active IOCs (Indicator of Compromise) for the Cyber Threat Intelligence Platform, where we can match the similar or same malicious incidents occurred to other customers based on the IOCs.

4. Evaluation

Holmes has been deployed in an enterprise environment, where it can read mirrored SMTP records from the Elastic-Search (ES) server. In this section, we first present some case studies on the malicious emails detected by Holmes and then show how correlation analysis can reveal the attack scenarios caused by malicious emails, and finally we compare Holmes with other popular commercial email detectors.

4.1. Case Studies

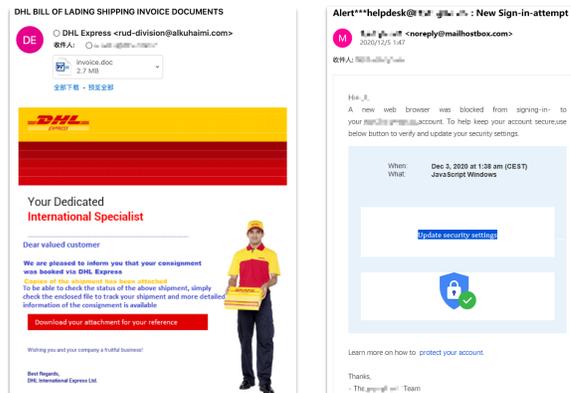
According to our monthly email system data, Holmes can discover around 1,000 anomalous emails each day. Among them, about 23% are truly malicious. And most of the malicious emails contain either phishing links or malware infected attachments. The rest are mainly spams and only a few are false positives.

Based on the detection results, we derive some malicious emails from our email server, which were not blocked by the anti-spam gateway but have been identified by our security analysts as the high risks, to reconstruct the attack stories. Here, we would showcase some delicate crafted phishing emails and describe their malicious behaviour in detail.

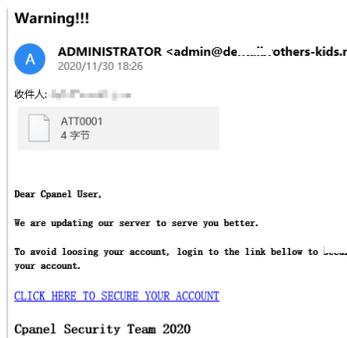
4.1.1. DHL Service. Fig 3 (a) shows a malicious email that pretends DHL service and plays following tricks: (1) The email uses a normally-seen subject that is associated with an invoice document; (2) The sender information has been modified as ‘DHL Express’, which can be implemented by some hacking tools, such as swaks [35] or cobalt strike [36]; (3) The email includes an attachment named *invoice.doc*, which is, in fact, a malicious Trojan document that utilizes the CVE-2017-11882 [37] vulnerability; (4) The email contains a delicate picture of DHL delivery service to spoof recipients.

In this attack scenario, an attacker who successfully exploited the vulnerability could run arbitrary code in the context of the current user (recipient). If the user is logged on with the administrative user rights, the attacker could take control of the affected system. The attacker could then install programs; view, change, or delete data; or create new accounts with full user rights. As we can see, users whose accounts are configured to have fewer user rights on the system could be less impacted than users who operate with administrative user rights.

4.1.2. New Sign-in Attempt. The email shown in Fig 3 (b) uses a deceptive subject named “New Sign-in Attempt”, aiming to spoof recipients to change their email account



(a) DHL Service (b) New Sign-in Attempt



(c) Warning !!

Figure 3. Phishing Email Samples

password. Once the recipient clicks the button of “Update security settings”, the web page will be redirected to the phishing website: [https://controladmin.7m.pl/login\[.\]html?#xxx@xxx.gov.xx](https://controladmin.7m.pl/login[.]html?#xxx@xxx.gov.xx), which induces the victim user to type in the username and password. The web page will, in the end, return to the enterprise homepage that the victim user works.

On the hacker side, the back-end server will receive the event log of the failed login attempts from the victim user, and then record the username and password. Hence, the hacker can use the legitimate email account to sign in, such as web page or email server, and can even further send an elaborately crafted phishing email to a person who is the victim’s frequent contact, which is hard to be detected by most security products.

4.1.3. Warning!!!. The email shown in Fig 3 (c) is similar to the attack shown in Fig 3(b) in that it also has a link embedded in the mail content for phishing campaign. However, the phishing link [https://armonaoil.com/admin/images/npgr/news/potcpanel\[.\]html](https://armonaoil.com/admin/images/npgr/news/potcpanel[.]html) is from a legitimate website rather than from a personally created malicious website, which indicates that the legitimate website has been compromised for the use of darknet market¹.

1. The darknet is most often used for illegal activities such as black markets, illegal file sharing, and the exchanging of illegal goods or services.

Email Subject	Microsoft	Tencent	Kaspersky	FireEye	McAfee	Qihoo-360	Holmes
SF Express New Order_INV 2019022411	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
Confirm Your Invoice for Payment	TRUE	FALSE	TRUE	TRUE	TRUE	FALSE	TRUE
Mail Update Notification. Inbox Full on	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Purchase Order	TRUE	FALSE	TRUE	TRUE	TRUE	FALSE	TRUE
About: Ownership Confirmation of***	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Reminder: Your Package Could Not Be Delivered***	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Re: **TOP URGENT** BL Draft Copy	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
REE:URGENT QUOTATION NEEDED ASAP	TRUE	FALSE	TRUE	TRUE	FALSE	FALSE	TRUE
Re:QUOTATION TEMPLATE2021	TRUE	FALSE	TRUE	TRUE	FALSE	FALSE	TRUE
***disconnected Fix Now!	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Validate your Password for***	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Your email***will be closed soon	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
Notifications undelivered emails to your mailbox	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE
DHL BILL OF LADING SHIPPING INVOICE	TRUE	FALSE	TRUE	TRUE	TRUE	TRUE	TRUE
Attention***mail upgrade	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	TRUE

Figure 4. Detection Result of Holmes Compared with Other Detectors in VirusTotal Enterprise Version

By further analysis, we found that the enterprise indeed opened the cPanel web hosting server for public access, which was vulnerable to the brute-force attack and remote external control. Furthermore, we examined the recent activities on some popular darknet markets and found that more than 3,000 sites of cPanel accesses were selling in the darknet market (raidforums) since 2020-11-22. Hence, it can be confirmed that the email attack is a phishing campaign caused by the third-party information leakage.

4.2. Attack Scenarios Revealed by Correlation Graph Analysis

As discussed in Section 3.4, email attacks can be clustered in line with some features, such as senders and subjects, and there may be geographical links with each other. In our investigation, we can clearly observe the email clusters on the directed-force graph generated by our CGA module, as has been demonstrated in Fig. 2. Those clusters can reveal some typical attack scenarios. Some are presented below.

4.2.1. Bitcoin Fraud. The attack scenario is related to a cluster which shows that hundreds of senders from different srcIp addresses sent the same email to spoof the recipients with an email content: “*Your computer has been controlled...transfer bitcoin to the wallet....*”.

Based on our experience of threat hunting, this situation indicates that a controlled botnet have been used for the email fraud.

4.2.2. Spams. The spams have been discovered in another cluster which shows that several srcIp addresses from a single country sent a large amount of spam emails that had similar interests of subjects, which involved with sensitive political topics, gamble and eroticism information. These emails are basically harmless but rather annoying users if the anti-spam gateway fails to filter them out.

However, it is still important to investigate those unusual email communication behaviours because some of the spam emails may involve with potential spy activities, which may lead to the information leakage of internal classified documents through social engineering.

4.2.3. Periodical Anomalous Behaviour. A notable clustered behaviour was also observed by the CGA module that a group of hackers from the same country using different srcIp addresses periodically sent malicious emails with similar subjects to our customers, where its phishing links use the same URL redirection technique. These emails use a crafted email content with the brand of targeted victims, such as the name of affiliations or IT support team. Once users follow their instructions to finish all actions, they will be redirected to their original organization’s homepage.

This periodical anomalous behaviour has eventually been identified as a long-term and targeted phishing activity by our security analysts.

4.3. A Comparative Study

To evaluate the detection capability of Holmes, we compare it with some commercial email detectors that are offered by six key security vendors in VirusTotal: Microsoft, Tencent, Kaspersky, FireEye, McAfee and Qihoo-360. We select 15 malicious emails as testing samples. They either contain a phishing link or have a malware infected attachment. All the testing samples were collected from the real-world threat hunting during the whole December month in 2020, and these samples had successfully bypassed the detection of the enterprise anti-spam gateway.

The comparison table is given in Fig 4, where the email subjects representing the 15 malicious emails are listed in the first column and the rest columns are the detection results from the commercial detectors and Holmes. A FLASE value

from a detector on a malicious email indicates that the detector failed to identify the malicious email.

From the table, we can see that Holmes can successfully detect all 15 malicious emails because the testing samples present strong behavioural deviations from the historical email records.

For the detectors of Microsoft, Kaspersky and FireEye, we can see a good performance on detection of those malicious emails that contain malware infected attachments. However, they fail to detect the malicious emails that contain phishing links. Based on the further analysis by our security experts, most of the phishing domains have been registered no more than three months and some of them are even from legitimate known enterprises. Moreover, all the phishing links include a specific URL to access the particular crafted phishing web page under the domain name that is shortly expired in around three days. Such a short-lived situation significantly increases the difficulty of anomaly detection.

From the comparison table, we can also see that McAfee demonstrates a moderate detection rate on the malicious emails that contain malware infected attachments, in which the No.8 and No.9 emails are failed to be detected. Similar to Microsoft, Kasperly and FireEye, McAfee also cannot detect the malicious emails that contain a newly registered phishing link.

Compared to all above detectors, Tencent and Qihoo-360 have a low detection rate. Among the 16 malicious emails, only two are detected by Tencent and three by Qihoo 360.

We would clarify that, the detection engines used for the comparison are supplied by the VirusTotal Enterprise Service. Since the version of the detectors may not be the same used in their commercial products, we would state that the comparison result cannot completely indicate the detection capability of their latest versions in the commercial products.

5. Enhancements in the Latest Implementation

After the first deployment to the enterprise environment, as mentioned in the above section, Holmes has been upgraded with a few enhancements.

5.1. Elastic-Search → Kafka Platform

In the latest version of Holmes, we rebuild the code warehouse that makes Holmes more efficient to discover anomalies in a much smaller rolling time window. The improvement is achieved by moving the data query system from Elastic-Search (ES) server to the real-time Kafka computing platform. The main difference between ES and Kafka is the way the data is processed. ES uses batch processing whereas Kafka uses stream processing, and the stream processing is more timely and efficient.

Kafka is an open-source distributed event streaming platform. It consists of producers, cluster (brokers) and consumers. Due to its high throughput and availability, Kafka has been widely used by thousands of companies for high-performance data pipelines, streaming analytics, data

integration, and mission-critical applications. In our case, the producers of Kafka are probers deployed in the enterprise network to sniff network traffic for the use of detection, brokers are the middle-ware mechanism to distribute data streams, and consumers are the detector of Holmes.

The advantage of using Kafka is that Holmes can detect anomalies in five minutes without the risk of server crash, significantly reducing the computing consumption. Furthermore, use of Kafka can also help our security analysts to better schedule time for threat identification and improve the efficiency of threat responses.

5.2. Email Direction Classification

Based on the report from our automated security operation center (SoC), most of the email threats are from the external sources, particularly from overseas, and the beacon of lateral phishing occurs with a much lower frequency than the inbound email threats, whereas the cases of data information leakage occurs with a medium frequency.

To improve the result interpret-ability, in the new implementation, we further classify the anomalous emails into three categories: inbound email threats, lateral phishing threats and outbound data leakage. The classification can help our security analysts to easily locate the potential threats so that the procedure of threat identification can be further accelerated.

With the new Holmes implementation, we are able to provide the daily threat report and offer the IOCs as a Software as a Service (SaaS) in a cloud platform to alert our customers so that they can response the threats in real time.

6. Related Work and Future Direction

Email is the most vulnerable entrance that can be easily utilized by the hacker for an attack, which, however, has not drawn sufficient attentions in many companies and organizations. In this section, we would briefly introduce some commercial products that particularly focus on email security and discuss some ideas for anomalous email detection, which can be the direction of future work.

6.1. Trustworthy Social Network Graph

User entity behaviour analytic (UEBA) is an essential method to effectively discover the anomalous activities that deviate the normal baseline. As we stated earlier, a stable sender to recipient relation often exists in the daily communication. Therefore, a trustworthy social network graph can be used to describe the confidence coefficient within the email communication. The social network graph can build a sender to recipient correlation based on their historical communication records; Once a new relation is discovered, the email can be marked as an abnormal behaviour for further investigation.

6.2. Content and Salutation Analysis

Content analysis is the easiest and direct method to identify whether an email is malicious. Some detectors may also analyze the name of salutation in the email to see whether it matches the recipient name to ensure it is not a mass mailing. By examining the keywords related to eroticism, gamble and money in email content, the analysis can effectively identify the malicious emails with a low FPR. However, many enterprises do not allow the third-part security product to access the email contents even if the contents are used for the genuine detection purpose, which significantly increases the cost and difficulties of detection. To address the problem, in recent years, many email detectors have considered to embed themselves within the email services of cloud vendors as the third-part extensions. The method can not only solve the conflicting results we often see between DR and FPR, but also save the cost of implementation.

6.3. Attachment Detection

Most of the email gateways have been equipped with the anti-virus software to prevent malware infected attachments. The anti-virus software can effectively prevent known malware but with false negatives for variations and some non-PE malware. In addition, hackers can also use many different or sophisticated methods to encrypt or obfuscate the malware, such as AES-128 or Base64, which allows the malware to escape detection from the anti-virus software. Furthermore, network traffic usually has a low visibility to some intrusions but the situation in endpoints is different. Once an attachment has been dropped or installed to the endpoint, we can clearly observe the execution of the commands that are from the malicious attachment. Hence, the linkage between network and endpoint is also useful to reveal some attack scenarios, which is worth to be further researched.

7. Conclusion

In this paper, we introduce Holmes, a lightweight semantic based anomalous email detector, which can effectively discover malicious emails in the real-world cyber threat hunting. Holmes also demonstrates a viable solution that successfully transfers AI technology to the cyber security field and makes an excellent trade-off between the cost of algorithmic consumption and the detection performance.

We measure the performance of Holmes, and compare its detection capability with several well-known commercial detectors offered by the security companies in VirusTotal. Our evaluation result shows that Holmes significantly outperforms those commercial products in all kinds of malicious attack scenarios, which demonstrates its practical values in the commercial competition.

8. Discussion

Part of the work has been successfully transferred and integrated into our security products and we are now happy

to open-source the code of prototype implementations. Since we published our evaluation results, we have received many questions regarding to our work. Here, we would like to share some FAQs and our replies from the discussions.

8.1. Q1: Do you only select the testing samples that may be in flavour your evaluation result?

Holmes is more like a hunting tool that can assist security experts to discover the most critical incidents based on a controlled number of alerts. The testing samples used are totally from our daily threat hunting in the wild. They were not purposely-made just in favor of our design. For comparison, we would also add that the tools offered by other security vendors run detection on the email body, whereas Holmes runs on the email header, which avoids the confidential issues related to accessing the email body. In addition, those tools aim to reduce FPR (hence, high FPR may be achieved at the cost of low DR). Holmes, on the other hand, aims to improve DR and may result in a high FPR.

8.2. Q2: Why do you not measure the overall FPR and TNR?

Following Question Q1, the purpose of developing Holmes is to discover those unknown threats that fail to be detected by the traditional email-body based tools . The main reason why we do not measure the FPR and TNR lies in the way of they are implemented. The tools mentioned in Fig 4 are deployed in the cloud, which requires to universally examine the malicious emails from a range of different sources. The diversity does not allow the detectors to generate too many false positives, leading to the poor DR on unknown threats. But Holmes is deployed locally in a specific enterprise environment, which enables its detection capability to be closely associated with the enterprise business. In addition, the number of emails received or sent daily in a single enterprise usually keeps a stable and controlled manner (around 5,000 to 10,000 based on the scale of different enterprise), in which situation the FPR and TNR are not the important metrics.

8.3. Q3: How about the comparison result of Holmes with the state-of-the-art in academia?

Abundant work can be found in the literature related to anomaly detection in the network security area. Anomalous email detection is now more targeted on industry application than research in academia, like the effort we made here. Therefore, we compare our design with the popular detection tools currently used.

Acknowledgement

The work is supported and funded by Sangfor Technologies Inc. and Ministry of Water Resources. Confidential information related to any personal and enterprise information

will not be available in public resources. The source code of Holmes can be used with the permission of the author (Peilun Wu).

No personally identifying information or sensitive data was shared with any non-employee of Sangfor Technologies Inc. and Ministry of Water Resources. Our project also received legal approval from Sangfor Technologies Inc. and Ministry of Water Resources, who had permissions to analyze and operate on the data. Our proposed anomalous email detector, Holmes, was deployed within the email security module of the cyber threat hunting platform, any detected attacks were reported to customers in real time to prevent further financial loss and harm.

References

- [1] I. D. Foster, J. Larson, M. Masich, A. C. Snoeren, S. Savage, and K. Levchenko, "Security by any other name: On the effectiveness of provider based email security," in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 450–464, 2015.
- [2] VirusTotal. "<https://www.virustotal.com/gui/home/search>". (accessed: 11.25.2020).
- [3] MITRE. "<https://attack.mitre.org/matrices/enterprise/>". (accessed: 11.25.2020).
- [4] G. Ho, A. Cidon, L. Gavish, M. Schweighauser, V. Paxson, S. Savage, G. M. Voelker, and D. Wagner, "Detecting and characterizing lateral phishing at scale," in *28th {USENIX} Security Symposium ({USENIX} Security 19)*, pp. 1273–1290, 2019.
- [5] T. Feng and C. Yue, "Visualizing and interpreting rnn models in url-based phishing detection," in *Proceedings of the 25th ACM Symposium on Access Control Models and Technologies*, pp. 13–24, 2020.
- [6] A. Blum, B. Wardman, T. Solorio, and G. Warner, "Lexical feature based phishing url detection using online learning," in *Proceedings of the 3rd ACM Workshop on Artificial Intelligence and Security*, pp. 54–60, 2010.
- [7] S. Garera, N. Provos, M. Chew, and A. D. Rubin, "A framework for detection and measurement of phishing attacks," in *Proceedings of the 2007 ACM workshop on Recurring malware*, pp. 1–8, 2007.
- [8] C. Whittaker, B. Ryner, and M. Nazif, "Large-scale automatic classification of phishing pages," 2010.
- [9] Barracuda. "<https://www.barracuda.com/glossary/email-spoofing>". (accessed: 11.25.2020).
- [10] M. Wong and W. Schlitt, "Sender policy framework (spf) for authorizing use of domains in e-mail, version 1," tech. rep., RFC 4408, April, 2006.
- [11] E. Allman, J. Callas, M. Delany, M. Libbey, J. Fenton, and M. Thomas, "Domainkeys identified mail (dkim) signatures," tech. rep., RFC 4871, May, 2007.
- [12] M. Kucherawy and E. Zwicky, "Domain-based message authentication, reporting, and conformance (dmarc)," ser. *RFC7489*, 2015.
- [13] J. Chen, V. Paxson, and J. Jiang, "Composition kills: A case study of email sender authentication," in *29th {USENIX} Security Symposium ({USENIX} Security 20)*, 2020.
- [14] P. Wu, N. Moustafa, S. Yang, and H. Guo, "Densely connected residual network for attack recognition," *19th IEEE International Conference on Trust, Security and Privacy in Computing and Communications (IEEE TrustCom)*, 2020.
- [15] S. Yang, P. Wu, and H. Guo, "Dualnet: Locate then detect effective payload with deep attention network," *IEEE Conference on Dependable and Secure Computing (DSC)*, 2021.
- [16] P. Wu, H. Guo, and N. Moustafa, "Pelican: A deep residual network for network intrusion detection," in *50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, pp. 55–62, IEEE, 2020.
- [17] P. Wu and H. Guo, "Lunet: A deep neural network for network intrusion detection," in *IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 617–624, IEEE, 2019.
- [18] I. Fette, N. Sadeh, and A. Tomasic, "Learning to detect phishing emails," in *Proceedings of the 16th international conference on World Wide Web*, pp. 649–656, 2007.
- [19] M. I. Jordan and T. M. Mitchell, "Machine learning: Trends, perspectives, and prospects," *Science*, vol. 349, no. 6245, pp. 255–260, 2015.
- [20] T. Tulabandhula and C. Rudin, "Machine learning with operational costs," 2013.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [22] W. U. Hassan, A. Bates, and D. Marino, "Tactical provenance analysis for endpoint detection and response systems," in *2020 IEEE Symposium on Security and Privacy (SP)*, pp. 1172–1189, IEEE, 2020.
- [23] X. Han, T. Pasquier, A. Bates, J. Mickens, and M. Seltzer, "Unicorn: Runtime provenance-based detector for advanced persistent threats," *arXiv preprint arXiv:2001.01525*, 2020.
- [24] S. Ma, X. Zhang, and D. Xu, "Protracer: Towards practical provenance tracing by alternating between logging and tainting.," in *NDSS*, 2016.
- [25] X. Han, T. Pasquier, and M. Seltzer, "Provenance-based intrusion detection: opportunities and challenges," in *10th {USENIX} Workshop on the Theory and Practice of Provenance (TaPP 2018)*, 2018.
- [26] I. U. Haq, I. Gondal, P. Vamplew, and S. Brown, "Categorical features transformation with compact one-hot encoder for fraud detection in distributed environment," in *Australasian Conference on Data Mining*, pp. 69–80, Springer, 2018.
- [27] G. Andresini, A. Appice, N. Di Mauro, C. Loglisci, and D. Malerba, "Exploiting the auto-encoder residual error for intrusion detection," in *2019 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pp. 281–290, IEEE, 2019.
- [28] V. P. KS and K. Gurumurthy, "Design of high performance quaternary adders," in *2011 41st IEEE International Symposium on Multiple-valued logic*, pp. 22–26, IEEE, 2011.
- [29] P. Cerda, G. Varoquaux, and B. Kégl, "Similarity encoding for learning with dirty categorical variables," *Machine Learning*, vol. 107, no. 8, pp. 1477–1494, 2018.
- [30] Y. Zhang, R. Jin, and Z.-H. Zhou, "Understanding bag-of-words model: a statistical framework," *International Journal of Machine Learning and Cybernetics*, vol. 1, no. 1-4, pp. 43–52, 2010.
- [31] Scikit-Learn. "<https://scikit-learn.org/>". (accessed: 11.25.2020).
- [32] Q. Le and T. Mikolov, "Distributed representations of sentences and documents," in *International conference on machine learning*, pp. 1188–1196, 2014.
- [33] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek, "Loop: local outlier probabilities," in *Proceedings of the 18th ACM conference on Information and knowledge management*, pp. 1649–1652, 2009.
- [34] T. M. Fruchterman and E. M. Reingold, "Graph drawing by force-directed placement," *Software: Practice and experience*, vol. 21, no. 11, pp. 1129–1164, 1991.
- [35] Swaks. "<https://github.com/jetmore/swaks>". (accessed: 11.25.2020).
- [36] CobaltStrike. "<https://www.cobaltstrike.com/>". (accessed: 11.25.2020).
- [37] CVE-2017-11882. "<https://msrc.microsoft.com/update-guide/en-US/vulnerability/CVE-2017-11882>". (accessed: 11.25.2020).