

# Understanding Attack Trends from Security Blog Posts Using Guided-topic Model

TATSUYA NAGAI<sup>1</sup> MAKOTO TAKITA<sup>1,†1</sup> KEISUKE FURUMOTO<sup>2</sup> YOSHIAKI SHIRAISHI<sup>1,a)</sup> KELIN XIA<sup>3</sup>  
 YASUHIRO TAKANO<sup>1</sup> MASAMI MOHRI<sup>4</sup> MASAKATU MORII<sup>1</sup>

Received: March 12, 2019, Accepted: September 11, 2019

**Abstract:** Organizations are plagued by sophisticated and diversified cyber attacks. In order to prevent such attacks, it is necessary to understand threat trends and to take measures to protect their assets. Security vendors publish reports which contain threat trends or analysis of malware. These reports are useful for help in responding to a cyber security incident. However, it is difficult to collect threat information from multiple sources such as security blog posts. In this paper, we propose a method to efficiently collect information from the relationship between words using SeededLDA. In our case studies, we visualize the relationship between the words from security blog posts which were published in 2017 by eight security vendors, and demonstrate how our method helps to understand threat trends in the IoT industry and financial institutions.

**Keywords:** security blog post, topic model, threat analysis, incident handling

## 1. Introduction

Organizations are plagued by sophisticated and diversified cyber-attacks. It is difficult to prevent such attacks only from existing systems such as Intrusion Detection System (IDS) and Security Information and Event Management (SIEM). In incident response, a security analyst quickly identifies infected systems and investigates the damage and plans a defense strategy. Not only that, but it is also necessary to be aware of threat trends and take proactive measures to protect their assets in order to respond quickly.

Threat information can be collected from such as security blog posts published by security vendors. In some cases, it is possible to obtain useful information from individuals with blogs or SNS. However, it is difficult to integrate threat information from multiple sources which have different labeling.

Security blog posts have become a vital source for input of existing systems. Liao et al. [1] proposed a method for extracting indicators such as IP addresses of malware and C2 servers from security blog posts and converting them into the openIOC format. Sabotke et al. [2] proposed a system that detects software vulnerabilities from Twitter streams utilizing existing vulnerability databases such as the National Vulnerability Database (NVD) and the Open Sourced Vulnerability Database (OSVDB). Sapienza et al. [3] proposed a system that analyzes online social media such as Twitter and dark web forums to generate alerts.

However, some attacks like Advanced Persistent Threat (APT) cannot be prevented by the existing systems. So, it is essential to manually integrate threat information from the security blog posts, and to understand threat trends and attack intent.

The topic model has been used to retrieve information from documents. Ma et al. [4] proposed a three-stage clustering method using Latent Dirichlet Allocation (LDA) [5] which is one of the topic models. This method finds an optimal number of topics to cluster documents. Onan et al. [6] proposed a K-means clustering using LDA. Tagarelli et al. [7] proposed a method to cluster text segments from multi-topic documents. The topic model allows us to analyze the topics which are contained the documents. However, the topic model might not generate the topic that a security analyst desired.

The purpose of research is establishing a method to predict threat trends from the security blog posts. In this paper, we propose a method to efficiently collect information by visualizing the relationship between words using SeededLDA [8]. In our case studies, we visualize the relationship between the words from the security blog posts which were published in 2017 by eight security vendors and show that our method helps to understand threat trends in the IoT industry and financial institutions.

## 2. Related Works

Bridges et al. [9] proposed a method to automatically label unstructured text such as NVD text description associated with structured data. Mulwad et al. [10] proposed a method to extract vulnerabilities and attacks using Wikitology. McNeil et al. [11] proposed a method to extract software names, vulnerability categories, vulnerability effects, and exploit effects. Joshi et al. [12] proposed a method to generate a Resource Description Framework (RDF) data, which is automatically processed, from un-

<sup>1</sup> Kobe University, Kobe Hyogo 657–8501, Japan

<sup>2</sup> National Institute of Information and Communications Technology, Koganei Tokyo 184–8795, Japan

<sup>3</sup> Nanyang Technological University, 637371 Singapore

<sup>4</sup> Gifu University, Gifu 501–1193, Japan

<sup>†1</sup> Presently with University of Hyogo

<sup>a)</sup> zenmei@port.kobe-u.ac.jp

structured data. Mittal et al. [13] proposed a method to alert for security analysts from the Twitter stream based on an organization’s system profile.

These methods label the threat information which contained in documents and allow security analysts to understand where the threat information is written.

However, they focus on single attack activity like malware, vulnerability, affected software, etc. To counter a series of cyber attacks, it is necessary to understand attack trends in time series data and the depth of the relationship between each attack by combining multiple security blog posts.

### 3. Topic Model and Guided-topic Model

A topic model is a generative model that deals with a document as a set of words, and a word is generated from topics that are clusters of semantically similar words.

**Figure 1** (a) shows graphical notation of LDA, one of the topic models. Let  $M$  be the number of documents,  $d$  be the document,  $T$  be the number of topics, and  $w = \{w_{d_1}, w_{d_2}, \dots, w_{d_n}\}$  be the words appearing in document  $d$ . The observed variable is colored grey in Fig. 1. LDA assumes the following generative process for each document  $d$ :

1. For each topic  $k = 1, 2, \dots, T$ ,
  - A) Choose a word distribution  $\varphi_k \sim \text{Dir}(\beta)$ .
2. For each document  $d$ ,
  - A) Choose a topic distribution  $\theta_d \sim \text{Dir}(\alpha)$ .
  - B) For each word  $w = \{w_{d_1}, w_{d_2}, \dots, w_{d_n}\}$ ,
    - i. Select a topic  $z_{d_i} \sim \text{Mult}(\theta_d)$ .
    - ii. Select a word  $w_{d_i} \sim \text{Mult}(\varphi_{z_{d_i}})$ .

The topic model associates the words with the topic by the word correlations. However, the topic generated by LDA may not be meaningful to a user. Jagarlamudi et al. proposed SeededLDA [8] that allows a user to give additional information to the topic model in order to learn topics of specific interest to a user. In SeededLDA, a user provides seed sets to guide the topics.

SeededLDA’s graphical notation is shown in Fig. 1 (b). SeededLDA assumes the following generative process for each document  $d$ :

1. For each topic  $k = 1, 2, \dots, T$ ,
  - A) Choose a regular word distribution  $\varphi_k^r \sim \text{Dir}(\beta^r)$ .
  - B) Choose a seed word distribution  $\varphi_k^s \sim \text{Dir}(\beta^s)$ .
2. For each seed set  $s = 1, 2, \dots, S$ ,
  - A) Choose a group-topic distribution  $\psi_s \sim \text{Dir}(\alpha)$ .
3. For each document  $d$ ,
  - A) Choose a seed binary vector  $\vec{b}$  of length  $S$ .
  - B) Choose a document-group distribution  $\zeta^d \sim \text{Dir}(\tau \vec{b})$ .
  - C) Choose a document-group variable  $g \sim \text{Mult}(\zeta^d)$ .
  - D) Choose a topic distribution  $\theta_d \sim \text{Dir}(\psi_g)$ .
  - E) For each word  $w = \{w_{d_1}, w_{d_2}, \dots, w_{d_n}\}$ ,
    - i. Select a topic  $z_{d_i} \sim \text{Dir}(\theta_d)$ .
    - ii. Select an indicator  $x_i \sim \text{Bern}(\pi_{z_i})$ .
    - iii. If  $x_i$  is 0, select a word  $w_{d_i} \sim \text{Mult}(\varphi_{z_{d_i}}^r)$ .
    - iv. If  $x_i$  is 1, select a word  $w_{d_i} \sim \text{Mult}(\varphi_{z_{d_i}}^s)$ .

**Table 1** shows an example of a seed sets. For example, consider a short document “IoT system deploy in Factory”. According to the seed sets from Table 1, we define a binary vec-

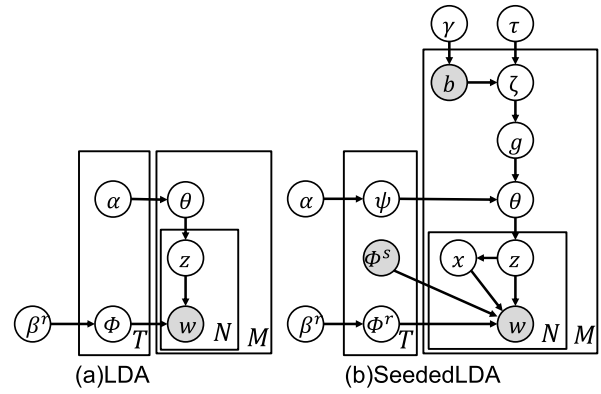


Fig. 1 Graphical notation of LDA and SeededLDA.

Table 1 Example of input of seed sets in SeededLDA.

1	mobile	mobile, smartphone, android, iphone
2	web server	web application, CMS, apache strut
3	IoT	wireless, IoT, connected device
4	medical	hospital, medical, health care, patient
5	factory	industrial control system, factory, plant

tor  $\vec{b} = \langle 0, 0, 1, 0, 1 \rangle$  because the document contains “IoT” whose topic number is 3 and “factory” whose topic number is 5. Seed vector  $\vec{b}$  is used to determine a document-group  $g$ . Documents which have the same vector are guided to be assigned the same topic. In addition, the seed word distribution  $\varphi_k^s$  restricts the model to select from the seed word set in the generative process. A probability of choosing the seed word distribution can be controlled by an indicator  $\pi$ . SeededLDA can be expected to classify documents as intended by a user without preparing a supervised data by providing the seed set.

### 4. Proposed Method: Understanding Attack Trends using SeededLDA

The purpose of research is establishing a method to analyze threat trends from the security blog posts. We propose a method to efficiently collect information by visualizing the relationship between words using SeededLDA.

Overview of the method is shown in Fig. 2. Inputs of the method are security blog posts and seed sets which an analyst provides. SeededLDA learns topics from the word correlations and the seed sets. Then, the word-topic distribution obtained by the SeededLDA with  $T$ -dimension is reduced to two-dimensional vectors and is visualized. Words close to each other in two-dimensional space are similar in terms of cyberspace. Our method has two approaches. One approach is to input all security blog posts in order to help an analyst understand the relationship between a specific industry and cyber-attacks such as attack methods and attack tools. Another approach is to input the security blog posts of a specific period in order to help an analyst understand threat trends in a specific industry.

#### 4.1 Preprocessing

Security vendors often advertise their products or themselves in security blog post and the advertising texts are not directly related to the threat information, so these posts are not covered in

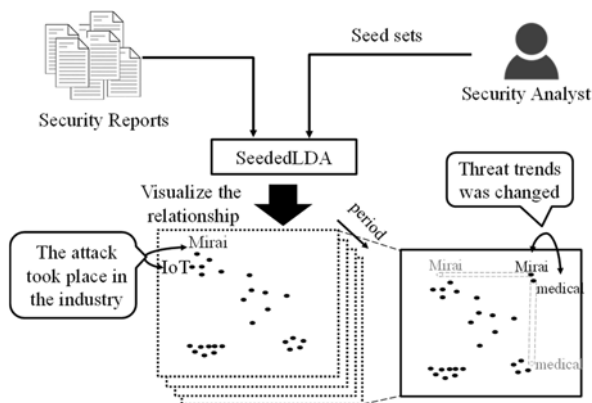


Fig. 2 Overview of proposed method.

Table 2 The list of additional stop-words.

APSB	XGen	cisco
CVE	ZDI	stealthwatch
deepsight	barracuda	symantec
Deep Security	blue coat	talos
Deep Discovery	bluecoat	next-gen
Folder Shield	TrendMicro	trend micro
InterScan Web Security	tushar richabadas	Control Manager
cloud generation firewall	Data Loss Prevention Manager	smart protection network
TM		

this paper. They also publish regular reports summarizing multiple vulnerabilities. Because it is difficult to extract detailed threat information about a single topic from these articles, these articles are also not covered in this paper.

The preprocessing is performed in the following procedures:

1. Enter a title, a body, and some figures of security blog posts.
2. Remove stop-words and special characters.
3. Extract compounds word using N-gram.

In this paper, step 2 and step 3 are executed by “Preprocess Text” module and “Extract N-Gram Features from Text” module on Microsoft Azure. In step 1, the titles of post and figures contain important information about the threat and should be entered together with the body. In step 2, stop-words, numbers, special characters, duplicate characters, email addresses, and URLs are removed from the inputs. At same time, the words in **Table 2** are entered into module as stop-words in the security field and these words are removed from the inputs. In addition, lemmatization, detection of sentence, and normalization to lowercase are performed. In step 3, N-gram is executed with  $n = 2$  to extract technical terms frequently displayed in security blog posts.

#### 4.2 SeededLDA Learning

We input the security blog posts and the seed sets to SeededLDA. An analyst can set the seed sets according to the state and an environment of the organization. Since the compound word is extracted by N-gram in Section 4.1, a compound word which contained the seed word is also set as the seed.

There may be a large number of security blog posts compared with a topic which we want to guide. In this case, we set multiple

seed topics with the same seed set. On the other hand, the number of seed words may be too large due to the size of the input documents and the setting of the seed sets. Because many seed words may reduce classification quality, we limit the number of seed words in each seed topic to a few dozen by sampling.

#### 4.3 Visualizing the Relationship between Words

We get the word-topic distributions from the learned SeededLDA in Section 4.2. The word-topic distribution is a  $T$ -dimension vector, in which  $T$  is the number of topics, so it is difficult to understand a relationship between words. Therefore, we reduce the word-topic distribution to two dimensions using t-SNE method [25] and visualize it to help an analyst understand cyber-attacks and threat trends. t-SNE can be used to compress data while preserving similarity between high-dimensional data and to plot closely related words on a two-dimensional plane.

##### 4.3.1 Understanding Attack Methods and Attack Tools

SeededLDA learns all security blog posts and the word-topic distribution is visualized in two dimensions.

An analyst focuses on a word which represents a specific industry and looks up a word which represents an attack method or a tool around it. If the word is found it indicates that an attack took place in that industry, etc.

##### 4.3.2 Understanding Threat Trends

SeededLDA learns security blog posts of a specific period and the word-topic distribution is visualized in two dimensions.

In the specific period, an analyst can find the relationship between attack methods and a specific industry. If the relationship is not found in another period, it can be understood that the attack was prevalent in the specific period and that the threat trends has changed.

### 5. Case Studies

#### 5.1 Setup

We performed case studies on whether our method provides information which helps an analyst understand attack methods and threat trends.

We apply the method described in Section 4.3.1 from a perspective of the IoT industry and a financial institution. We also apply the method described in Section 4.3.2 from a perspective of a financial institution. According to an annual report published by Symantec on Ref. [14], they found a 600 percent increase in overall IoT attacks in 2017. Because an analyst may want to know what attacks occur in the IoT industry, we apply the method to the IoT industry. They also reported two financial trojans are taken down towards the end of 2016. Because we expect to be able to observe the change of threat trends in a financial institution, we apply the method to a financial institution. In applying the method described in Section 4.3.2, the threat trend is observed by changing the input of the SeededLDA, where a window size of the input is three months and the shift width is one week.

Dataset is a collection of 875 reports published by TrendMicro [15], Cisco [16], Symantec [17], Barracuda [18], Druva [19], FireEye [20], Arbor [21], and Palo Alto [22] in 2017. In our case studies, we use the NLTK [23], a natural language toolkit, for preprocessing. GuidedLDA library [24] is used to perform the

**Table 3** Seed sets to input SeededLDA in case studies.

Guide topics	Seed set
DoS	denial of service, dos
Data breach	data breach
Infect	infect
Financial loss	money, financial
Credential	credential

SeededLDA. Seed sets to input SeededLDA are shown in **Table 3**. The number of topics  $T$  in SeededLDA are 36.  $\alpha$ ,  $\beta$ , and the probability of choosing the seed word distribution are set to 0.01, 0.01, and 0.7, respectively. We use the t-SNE method [25] to visualize the word-topic distribution. Scikit-learn library [26] is used to perform the t-SNE. Considering visibility, we only visualize words that appear 5–100 times in all documents.

**5.2 Understanding Attack Methods in the IoT Industry**

**Figure 3** shows a visualization of the word-topic distributions in two-dimension and an enlarged view of the word “internet thing” that is a word relating to the IoT industry. **Table 4** shows a list of terms contained in a topic related to “internet thing”. They are indexed in descending order of term frequency.

In Fig. 3, “router”, “printer”, “camera”, and “mirai” appear near “internet thing”. This is one of the IoT devices, so we confirmed that the relevant words appear nearby in visualization. The word, “firmware” also appears nearby. So, it can be estimated that some attack methods on IoT equipment use firmware.

According to a security blog post published by Cisco on April 18, 2017 [27], the practices of IoT security are summarized, including firmware updates. Therefore, it can be understood that there is an attack method aimed at the firmware on the IoT industry. In Fig. 3, the relationship between the IoT and the firmware appears, so we confirmed that our method provides information which helps an analyst understand the attack methods.

**5.3 Understanding Threat Trends in the IoT Industry**

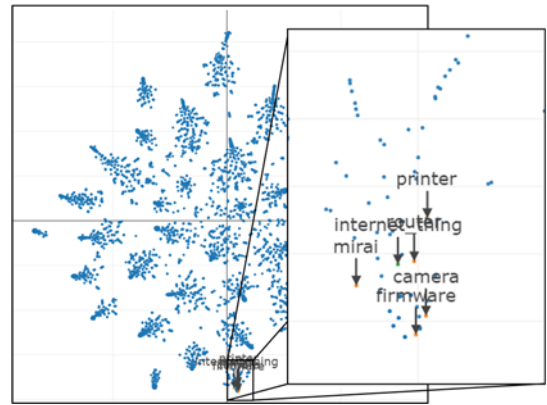
**Figures 4 and 5** show a visualization of the word-topic distributions in two-dimension during the period from June to September and the period from September to December they are an enlarged view of the word around “internet thing”.

Like Fig. 3, “router”, “printer”, “camera”, “firmware”, and “mirai” appear near “internet thing” in Figs. 4 and 5. Mirai which is a malware targeted at IoT devices spread worldwide in 2017. Therefore, IoT devices such as “router” and “printer” that were targeted by Mirai appear near “internet thing” throughout the year.

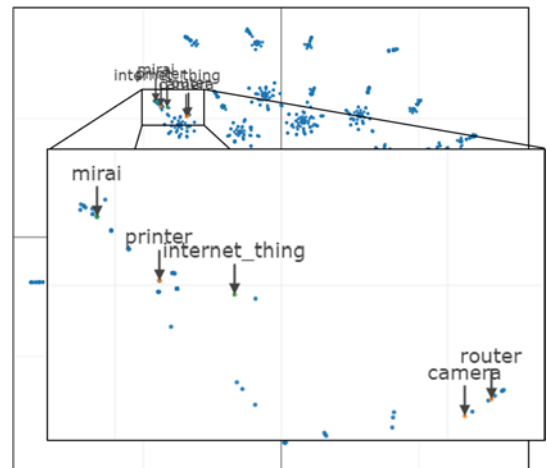
**5.4 Understanding Attack Methods in a Financial Institution**

**Figure 6** shows a visualization of the word-topic distributions in two-dimension and an enlarged view of the word “financial institution”. **Table 5** shows a list of terms contained in a topic related to “financial institution”. They are indexed in descending order of term frequency.

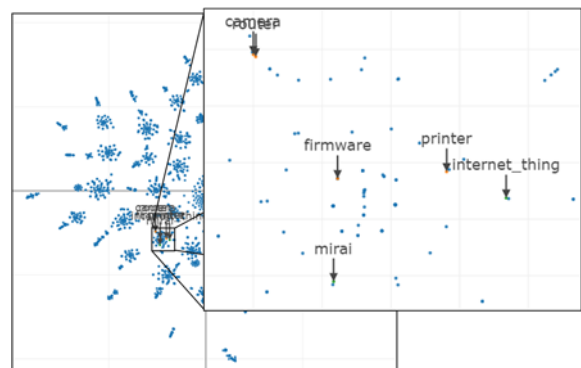
In Fig. 6, “bank account”, “steal malware”, “personal information”, and “fraudster” appear near “financial institution”. So, it



**Fig. 3** The relationship between words about the IoT industry.



**Fig. 4** The relationship between words about IoT industry during the period from June to September.



**Fig. 5** The relationship between words about IoT industry during the period from September to December.

can be estimated that a theft attack of a personal information is done to the financial institution.

According to a security blog post published by Cisco on October 11, 2017 [28], the importance of a stateful firewall is described. The security blog post says one of the reasons is that stealthy malware theft credit card numbers or financial institution passwords. So, we confirmed that our method provides information which helps the analyst understand the attack methods.

On the other hand, a security blog post published by Symantec on June 1, 2017 [29] summarizes malware targeting financial institutions. The report says that there was a lot of detection of Trojan horse *Ramnit* aiming at online banking. Unfortunately, we

**Table 4** Terms contained in a topic related to *internet thing*.

1	ascii_wide	26	edge_technology	51	connect_device	76	manufacturer
2	hajime	27	time_device	52	thing_internet	77	wide
3	default_password	28	webcam	53	video_camera	78	printer
4	device_internet	29	plant_floor	54	device_level	79	device
5	firmware	30	u.s._city	55	device_right	80	controller
6	casa	31	insecure_device	56	database_internet	81	transportation
7	industrial_system	32	manufacture_plant	57	system_critical	82	internet
8	mirai_botnet	33	step_risk	58	organization_device	83	manufacture
9	zombie	34	expose_system	59	stephen	84	home
10	when_internet	35	device_default	60	connect_thing	85	connect
11	password_device	36	because_device	61	computer_system	86	password
12	pipe	37	strong_unique	62	municipality	87	shutdown
13	urban	38	cable	63	router	88	vulnerable_device
14	camera	39	child_personal	64	internet_thing	89	equipment
15	device_secure	40	nonprofit_mission	65	city	90	sensor
16	thermostat	41	cyber_asset	66	home_device	91	wireless
17	unique_password	42	dvrs	67	industrial	92	thing
18	afterthought	43	program_nonprofit	68	mirai	93	critical_infrastructure
19	most_expose	44	cyber_vital	69	inappropriate	94	default
20	home_router	45	audio	70	company_device	95	system
21	specific_protocol	46	posture_device	71	secure_firmware	96	york_york
22	privacy_setting	47	urban_area	72	expose	97	physical
23	iiot	48	insecurity	73	ascii	98	vulnerable
24	inappropriate_content	49	world_wide	74	mongodb	99	insecure
25	shodan	50	device_year	75	plant	100	minimum

**Table 5** Terms contained in a topic related to *financial institution*.

1	holiday	26	underground_market	51	retailer	76	Gift
2	reputable	27	website_user	52	underground	77	Scam
3	cybercriminal_underground	28	newspaper	53	that_hacker	78	crime
4	mario	29	tempt	54	netizens	79	deal
5	counterfeit	30	threat_cyber	55	reputable_provider	80	cyber_crime
6	unsolicited	31	attack_long	56	possible_threat	81	fake
7	holiday_season	32	unsolicited_email	57	shop	82	prediction
8	attachment_link	33	that_criminal	58	scammer	83	unsuspecting
9	bank_account	34	cybercriminal_activity	59	season	84	kind
10	gift_card	35	honeypot	60	fraudulent	85	consumer
11	questionable	36	style_attack	61	financial_inst itution	86	dark
12	black_market	37	cryptolocker	62	fraudster	87	credit
13	expert_today	38	criminal_gang	63	bank	88	popular
14	police_department	39	however_threat	64	hacker	89	credit_card
15	online_shop	40	that_site	65	criminal	90	personal_information
16	online_account	41	heist	66	cybercrime	91	personal
17	open_attachment	42	malicious_purpose	67	west	92	news
18	unsuspecting_user	43	fleming	68	cybercriminal	93	online
19	north_american	44	steal_malware	69	cyber_criminal	94	fraud
20	holiday_shop	45	hacker_attack	70	internet_user	95	percent
21	online_bank	46	personal_financial	71	money	96	black
22	cyber_monday	47	employee_device	72	lucrative	97	steal
23	kind_attack	48	public_wifi	73	financial	98	account
24	since_then	49	possession	74	chat	99	transaction
25	west_african	50	wide_open	75	card	100	social

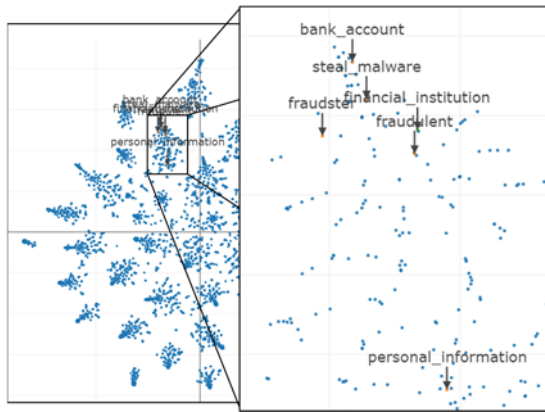


Fig. 6 The relationship between words about financial institution.

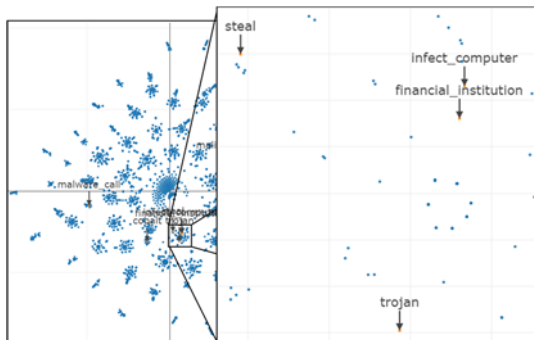


Fig. 7 The relationship between words about financial institution during the period from January to March.

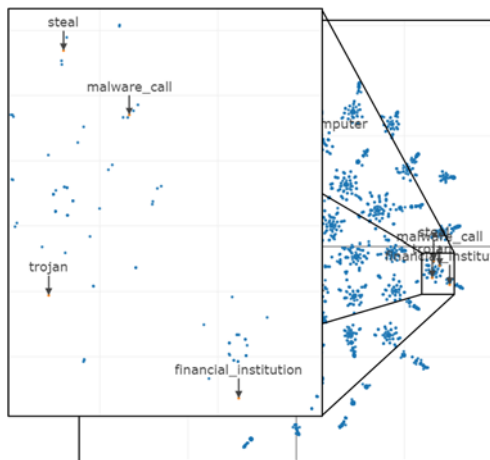


Fig. 8 The relationship between words about financial institution during the period from March to May.

confirmed that “trojan” which represents a Trojan horse did not appear near “financial institution” in Fig. 6. A consideration on this cause is given in Section 6.

### 5.5 Understanding Threat Trends in a Financial Institution

We observed the threat trends by changing the input of the SeededLDA, where a window size of input is three months and the shift width is one week.

Figures 7, 8, and 9 show a visualization of the word-topic distributions in two-dimension during the period from January to March, March to May, and September to November and they are an enlarged view of the word “financial institution”.

In Figs. 7 and 8, between January and May, the words “infect

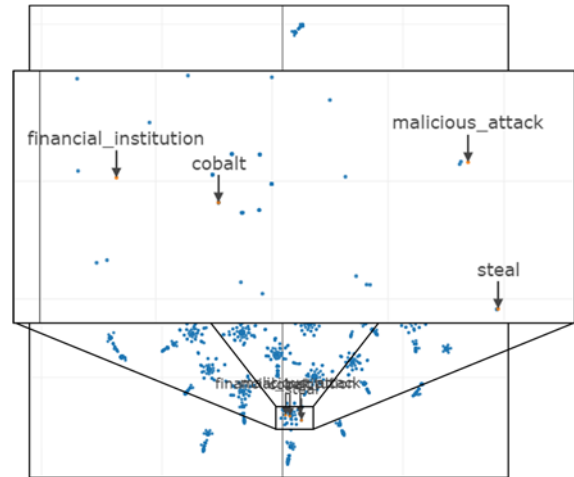


Fig. 9 The relationship between words about financial institution during the period from September to November.

computer” which is related to malware appear nearby “financial institution”. “trojan” also appears nearby it. While these words appear for most of the periods, in Fig. 9, we confirmed that the word “cobalt” appears nearby it.

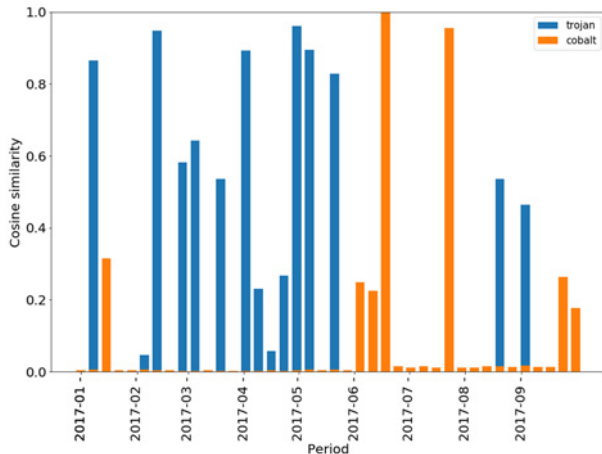
A security blog post published by TrendMicro on September 26, 2017 [30] says that cybercriminals used Cobalt Strike, which is a penetration testing tool, to attack ATMs. Therefore, we can get information about the newly emerged attack by visualizing the relationship between the words in each period. By comparison of the visualization in each period, we can also get information about when the attack was prevalent.

## 6. Discussion

In our case studies, we confirmed that our method provides information to help an analyst understand the attack methods and the threat trends by visualizing the relationship between the words. In the first case, the terms related to IoT devices appear near “internet thing” throughout the year because there were multiple posts about “Mirai” and its variants during 2017. In the second case, the relationship between “malware” and “trojan” which widely spread was not properly obtained. In order to analyze the cause, Fig. 10 shows the cosine similarity of the word-topic distributions of “financial institution”, “trojan” and “cobalt”, which were confirmed in Section 5.5. In Fig. 10, “trojan” is a high similarity between January and June. Although the similarity should be high in this period, there is the significant point that the similarity is 0 in the period. This is because the trojan horse attack was carried out in several industries, so it is thought that “financial institution” and “trojan” were not assigned to the same topic in SeededLDA learning phase. Our case studies did not capture this relationship well. If we set the seed word appropriately, these words will be assigned to the same topic.

Tables 4 and 5 show that there are several terms that are semantically similar to each other, but since they include unrelated terms, technical development is required to gather relevant terms in the security field.

In visualizing, some of the words have nothing to do with an attack method or a tool. In order to exclude these words in visualize, we can apply Named Entity Recognition (NER), which



**Fig. 10** The cosine similarity of the topic distributions of “financial institution”, “trojan” and “cobalt”.

is utilized in related studies [1], [9], [12], [13], to automatically extract words which are relevant to the threat information.

## 7. Conclusion

In this paper, we propose the method to efficiently gather information using SeededLDA. In our case studies, we confirm that our method provides information to help an analyst understand the attack methods and the threat trends by visualizing the relationship between the words.

Topics for future work include investigating the setting of seed sets and automatically finding words which are relevant to the threat information. Further studies are needed in order to automatically issue reports that the analyst needs to read.

## References

- [1] Liao, X., Yuan, K., Wang, X., Li, Z., Xing, L. and Beyah, R.: Acing the IOC game: Toward automatic discovery and analysis of open-source cyber threat intelligence, *Proc. 2016 ACM SIGSAC Conference on Computer and Communications Security*, pp.755–766 (2016).
- [2] Sabottke, C., Suci, O. and Dumitras, T.: Vulnerability Disclosure in the Age of Social Media: Exploiting Twitter for Predicting Real-World Exploits, *USENIX Security Symposium*, pp.1041–1056 (2015).
- [3] Sapienza, A., Bessi, A., Damodaran, S., Shakarian, P., Lerman, K. and Ferrara, E.: Early warnings of cyber threats in online discussions, *IEEE International Conference on Data Mining Series (ICDM)*, pp.667–674 (2017).
- [4] Ma, Y., Wang, Y. and Jin, B.: A three-phase approach to document clustering based on topic significance degree, *Journal of Expert Systems with Applications*, No.18, pp.8203–8210 (2014).
- [5] Blei, D.M., Ng, A.Y. and Jordan, M.I.: Latent Dirichlet Allocation, *Journal of Machine Learning Research*, No.3, pp.993–1022 (2003).
- [6] Onan, A., Bulut, H. and Korukoglu, S.: An improved ant algorithm with LDA-based representation for text document clustering, *Journal of Information Science*, Vol.43, No.2, pp.275–292 (2016).
- [7] Tagarelli, A. and Karypis, G.: A segment-based approach to clustering multi-topic documents, *Journal of Knowledge and Information Systems*, Vol.34, No.3, pp.563–595 (2013).
- [8] Jagarlamudi, J., Daumé III Hal. and Udupa, R.: Incorporating lexical priors into topic models, *Proc. 13th Conference of the European Chapter of the Association for Computational Linguistics*, pp.204–213 (2012).
- [9] Bridges, R.A., Jones, C.L., Iannacone, M.D. and Goodall, J.R.: Automatic Labeling for Entity Extraction in Cyber Security, arXiv preprint, arXiv:1308.4941 (2013).
- [10] Mulwad, V., Li, W., Joshi, A., Finin, T. and Viswanathan, K.: Extracting information about security vulnerabilities from web text, *Proc. 2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*, Vol.3, pp.257–260 (2011).
- [11] McNeil, N., Bridges, R.A., Iannacone, M.D., Czejo, B., Perez, N. and Goodall, J.R.: Pace: Pattern accurate computationally efficient

- bootstrapping for timely discovery of cyber-security concepts, *2013 12th International Conference on Machine Learning and Applications (ICMLA)*, Vol.2, pp.60–65 (2013).
- [12] Joshi, A., Lal, R., Finin, T. and Joshi, A.: Extracting cybersecurity related linked data from text, *IEEE 7th International Conference on Semantic Computing*, pp.252–259 (2013).
- [13] Mittal, S., Das, P.K., Mulwad, V., Joshi, A. and Finin, T.: Cybertwitter: Using twitter to generate alerts for cybersecurity threats and vulnerabilities, *Proc. 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pp.860–867 (2016).
- [14] Symantec: Internet Security Threat Report (ISTR) 2018, available from (<https://www.symantec.com/security-center/threat-report>) (accessed 2018-11-12).
- [15] TrendMicro: Simply Security News, Views and Opinions from Trend Micro, Inc., available from (<http://blog.trendmicro.com/>) (accessed 2018-10-30).
- [16] Cisco: Cisco Blog, available from (<https://blogs.cisco.com/>) (accessed 2018-10-20).
- [17] Symantec: Symantec Blogs, available from (<https://www.symantec.com/blogs/>) (accessed 2018-10-30).
- [18] Barracuda: Barracuda – Security, Access and Reliability for Cloud-Connected Networks and Applications, available from (<https://blog.barracuda.com/>) (accessed 2018-10-30).
- [19] Druva: Druva Blog: Data Protection and Beyond, available from (<https://www.druva.com/blog/>) (accessed 2018-10-30).
- [20] FireEye: Threat Research, available from (<https://www.fireeye.com/blog/threat-research.html>) (accessed 2018-10-30).
- [21] Arbor: Network Security Blog, available from (<https://asert.arbornetworks.com/>) (accessed 2018-10-30).
- [22] Palo Alto: Palo Alto Networks Blog, available from (<https://researchcenter.paloaltonetworks.com/>) (accessed 2018-10-30).
- [23] NLTK Project: Natural Language Toolkit — NLTK 3.3 documentation, available from (<http://www.nltk.org/>) (accessed 2018-10-30).
- [24] vi3k6i5/GuidedLDA: Semi supervised guided topic model with custom guidedLDA, available from (<https://github.com/vi3k6i5/GuidedLDA>) (accessed 2018-10-30).
- [25] Maaten, L.V.D. and Hinton, G.: Visualizing data using t-SNE, *Journal of machine learning research*, Vol.9, pp.2579–2605 (2008).
- [26] scikit-learn: machine learning in Python, available from (<http://scikit-learn.org/>) (accessed 2018-10-31).
- [27] Cisco: Demanding a Plan for Cyber Resilience in the IoT, available from (<https://blogs.cisco.com/security/demanding-a-plan-for-cyber-resilience-in-the-iot>) (accessed 2018-10-30).
- [28] Cisco: How is a Stateful Firewall like a Vintage Porsche?, available from (<https://blogs.cisco.com/security/how-is-a-stateful-firewall-like-a-vintage-porsche>) (accessed 2018-10-30).
- [29] Symantec: Financial malware more than twice as prevalent as ransomware, available from (<https://www.symantec.com/connect/blogs/financial-malware-more-than-twice-prevalent-ransomware>), accessed (accessed 2018-10-30).
- [30] TrendMicro: Attack The Machines: The lucrative business of ATM malware, available from (<https://blog.trendmicro.com/attack-machines-lucrative-business-atm-malware>) (accessed 2018-10-30).



**Tatsuya Nagai** received his B.E. and M.E. degrees from Kobe University, Japan, in 2017 and 2019, respectively.



**Makoto Takita** received his B.E., M.E., and D.E. degrees from Kobe University, Japan, in 2014, 2015, and 2018, respectively. He was a Researcher at the Graduate School of Engineering, Kobe University, Japan, in 2018. Since 2019, he has been an assistant professor at the School of Social Information Science, University

of Hyogo, Japan. His research interests include coding theory, information networks, and information security.



**Keisuke Furumoto** received his B.E., M.E., and D.E. degrees from Kobe University, Japan, in 2013, 2014 and 2018, respectively. He joined the National Institute of Information and Communications Technology (NICT) in 2018. His current research interests include information security and machine learning.



**Yoshiaki Shiraishi** received his B.E. and M.E. degrees from Ehime University, Japan, and the Ph.D. degree from the University of Tokushima, Japan, in 1995, 1997, and 2000, respectively. From 2002 to 2006 he was a lecturer at the Department of Informatics, Kindai University, Japan. From 2006 to 2013 he was an associate professor at the Department of Computer Science and Engineering, Nagoya Institute of Technology, Japan. Since 2013, he has been an associate professor at the Department of Electrical and Electronic Engineering, Kobe University, Japan. His current research interests include information security, cryptography, computer network, and knowledge sharing and creation support. He received the SCIS 20th Anniversary Award and the SCIS Paper Award from ISEC group of IEICE in 2003 and 2006, respectively. He received the SIG-ITS Excellent Paper Award from SIG-ITS of IPSJ in 2015. He is a member of IEEE, ACM, and a senior member of IPSJ and IEICE.

associate professor at the Department of Computer Science and Engineering, Nagoya Institute of Technology, Japan. Since 2013, he has been an associate professor at the Department of Electrical and Electronic Engineering, Kobe University, Japan. His current research interests include information security, cryptography, computer network, and knowledge sharing and creation support. He received the SCIS 20th Anniversary Award and the SCIS Paper Award from ISEC group of IEICE in 2003 and 2006, respectively. He received the SIG-ITS Excellent Paper Award from SIG-ITS of IPSJ in 2015. He is a member of IEEE, ACM, and a senior member of IPSJ and IEICE.



**Kelin Xia** received his Ph.D. degree from Graduate University of Chinese Academy of Sciences, China in 2013. He is currently with Nanyang Technological University as an assistant professor. His research interests include topological data analysis, mathematical modeling of biomolecular systems, and scientific computing.

computing.



**Yasuhiro Takano** received his Ph.D. (Info. Sc.) and Dr.Sc. (Tech.) degrees, respectively, from Japan Advanced Institute of Science and Technology (JAIST) and the University of Oulu, Finland, in 2016. He is currently with Kobe University as an assistant professor. His research interests include signal processing for communica-

tions engineering.



**Masami Mohri** received her B.E. and M.E. degrees from Ehime University, Japan, in 1993 and 1995 respectively. She received the Ph.D. degree in Engineering from the University of Tokushima, Japan in 2002. From 1995 to 1998 she was an assistant professor at the Department of Management and Information Science, Kagawa Junior College, Japan. From 1998 to 2002 she was a research associate of the Department of Information Science and Intelligent Systems, the University of Tokushima, Japan. From 2003 to 2007 she was a lecturer of the same department. From 2007 to 2017, she was an associate professor at the Information and Multimedia Center, Gifu University, Japan. Since 2017, she has been an associate professor at the Department of Electrical, Electronic and Computer Engineering in the same university. Her research interests are in coding theory, information security, and cryptography. She is a member of IEEE and a senior member of IEICE.

Kagawa Junior College, Japan. From 1998 to 2002 she was a research associate of the Department of Information Science and Intelligent Systems, the University of Tokushima, Japan. From 2003 to 2007 she was a lecturer of the same department. From 2007 to 2017, she was an associate professor at the Information and Multimedia Center, Gifu University, Japan. Since 2017, she has been an associate professor at the Department of Electrical, Electronic and Computer Engineering in the same university. Her research interests are in coding theory, information security, and cryptography. She is a member of IEEE and a senior member of IEICE.



**Masakatu Morii** received his B.E. degree in electrical engineering and the M.E. degree in electronics engineering from Saga University, Saga, Japan, and the D.E. degree in communication engineering from Osaka University, Osaka, Japan, in 1983, 1985, and 1989, respectively. From 1989 to 1990 he was an Instructor in the Department of Electronics and Information Science, Kyoto Institute of Technology, Japan. From 1990 to 1995 he was an Associate Professor at the Department of Computer Science, Faculty of Engineering, Ehime University, Japan. From 1995 to 2005 he was a Professor at the Department of Intelligent Systems and Information Science, Faculty of Engineering, the University of Tokushima, Japan. Since 2005, he has been a Professor at the Department of Electrical and Electronic Engineering, Faculty of Engineering, Kobe University, Japan. His research interests are in error correcting codes, cryptography, discrete mathematics, and computer networks and information security. He is a member of the IEEE and a fellow of IEICE.

From 1989 to 1990 he was an Instructor in the Department of Electronics and Information Science, Kyoto Institute of Technology, Japan. From 1990 to 1995 he was an Associate Professor at the Department of Computer Science, Faculty of Engineering, Ehime University, Japan. From 1995 to 2005 he was a Professor at the Department of Intelligent Systems and Information Science, Faculty of Engineering, the University of Tokushima, Japan. Since 2005, he has been a Professor at the Department of Electrical and Electronic Engineering, Faculty of Engineering, Kobe University, Japan. His research interests are in error correcting codes, cryptography, discrete mathematics, and computer networks and information security. He is a member of the IEEE and a fellow of IEICE.