



This is a repository copy of *Know abnormal, find evil: frequent pattern mining for ransomware threat hunting and intelligence*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/128370/>

Version: Accepted Version

Article:

Homayoun, S., Dehghantanha, A., Ahmadzadeh, M. et al. (2 more authors) (2020) Know abnormal, find evil: frequent pattern mining for ransomware threat hunting and intelligence. *IEEE Transactions on Emerging Topics in Computing*, 8 (2). pp. 341-351. ISSN 2168-6750

<https://doi.org/10.1109/TETC.2017.2756908>

© 2017 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other users, including reprinting/ republishing this material for advertising or promotional purposes, creating new collective works for resale or redistribution to servers or lists, or reuse of any copyrighted components of this work in other works. Reproduced in accordance with the publisher's self-archiving policy.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Know Abnormal, Find Evil: Frequent Pattern Mining for Ransomware Threat Hunting and Intelligence

Sajad Homayoun, Ali Dehghantanha, Marzieh Ahmadzadeh, Sattar Hashemi, Raouf Khayami

Abstract—Emergence of crypto-ransomware has significantly changed the cyber threat landscape. A crypto ransomware removes data custodian access by encrypting valuable data on victims' computers and requests a ransom payment to re-instantiate custodian access by decrypting data. Timely detection of ransomware very much depends on how quickly and accurately system logs can be mined to hunt abnormalities and stop the evil. In this paper we first setup an environment to collect activity logs of 517 *Locky* ransomware samples, 535 *Cerber* ransomware samples and 572 samples of *TeslaCrypt* ransomware. We utilize Sequential Pattern Mining to find Maximal Frequent Patterns (MFP) of activities within different ransomware families as candidate features for classification using J48, Random Forest, Bagging and MLP algorithms. We could achieve 99% accuracy in detecting ransomware instances from goodware samples and 96.5% accuracy in detecting family of a given ransomware sample. Our results indicate usefulness and practicality of applying pattern mining techniques in detection of good features for ransomware hunting. Moreover, we showed existence of distinctive frequent patterns within different ransomware families which can be used for identification of a ransomware sample family for building intelligence about threat actors and threat profile of a given target.

Index Terms—Malware, ransomware, crypto ransomware, ransomware detection, ransomware family detection.

I. INTRODUCTION

CYBERCRIMINALS pose a real and persistent threat to business, government and financial institutions all around the globe [1]. The volume, scope and cost of cybercrime all remain on an upward trend [2]. Malicious programs have always been an important tool in cyber criminals portfolios and almost everyday we are detecting new variants of malware programs [3]. Development and wide adoption of e-currencies such as Bitcoin led to many changes in cybercriminal activities including development of a new type of malware called ransomware [4]. Ransomware is a type of malware that removes a custodian access to her data and request for a ransom payment to re-instantiate data access [5]. There are two main types of ransomwares namely *Locker* and *Crypto* ransomwares. The former locks a system and denies users' access without making any changes to the data stored on the system while the latter encrypts all or selected data usually

using a strong cryptography algorithm such as AES or RSA [6].

Ransomware has dominated the threat landscape in 2016 with annual increase rate of 267% [7]. It is estimated that in 2014 only, cybercriminals have made more than \$3 million profit using ransomware programs [8]. These days, ransomware programs are indiscriminately targeting all industries ranging from healthcare to the banking sector and even power grids [2]. The Crypto-ransomware programs are much more popular than Lockers as almost always security engineers could find ways to unlock a system without paying the ransom while the only viable solution for decrypting strongly encrypted data is to pay ransom and receive decryption key [9]. Therefore, focus of this paper is only on crypto-ransomware and in the rest of the paper, the word "ransomware" is actually referring to the "crypto-ransomware" only. It was already reported that cyber security training and employee awareness would reduce the risk of ransomware attacks [10]. However, automated tools and techniques are required to detect ransomware applications before they are launched [11] or within a short period after their execution [12]. The growing danger of ransomware attacks requires new solutions for prevention, detection and removing ransomwares programs.

In this paper, we are using a sequential pattern mining technique to detect best features for classification of ransomware applications from benign apps as well as identifying a ransomware sample family. We investigate usefulness of our detected features by applying them in *J48*, *Random Forest*, *Bagging* and *MLP* classification algorithms against a dataset contains 517 *Locky* ransomware samples, 535 *Cerber* ransomware samples, 572 samples of *TeslaCrypt* ransomware and 220 standalone Windows Portable and Executable (PE32) benign applications. We not only achieved 99% accuracy in detection of ransomware samples and 96.5% in detection of their families but reduced the detection time to less than 10 seconds of launching a ransom application; a third of the time reported by earlier studies i.e. [13]. Our results are not only indicative of usefulness of pattern mining techniques in identification of best features for hunting ransomware applications but show how patterns of different ransomware families can help in detecting a ransomware family which assist in building intelligence about threats applicable to a given target. To the best of authors knowledge this is the very first paper applying sequence pattern mining to detect frequent features of ransomware applications and to build vectored datasets of ransomware applications logs. Our created datasets contain

S. Homayoun, M. Ahmadzadeh and Raouf Khayami are with the Department of IT and Computer Engineering, Shiraz University of Technology, Shiraz, Iran. e-mail: S.Homayoun@sutech.ac.ir.

A. Dehghantanha is with Department of Computer Science, School of Computing, Science and Engineering, University of Salford, Salford, U.K.

S. Hashemi is with Department of Computer Engineering, Shiraz University, Shiraz, Iran.

logs of Dynamic Link Libraries (DLL) activities, file system activities and registry activities of 1624 ransomware samples from three different families and 220 benign applications.

We are using widely accepted criteria namely True Positive (TP), False Positive (FP), True Negative (TN), and False Negative to evaluate our model [14]–[16]. TP is reflecting total samples that correctly identified. FP shows incorrectly identified samples. TN demonstrates the number of correctly rejected samples, while FN shows incorrectly rejected samples. *Precisions* of a classification algorithm is a measure of relevancy of results and is calculated by dividing TP by total of FP and TP predicted by a classifier as shown in equation (1). *Recall* reflects the proportion of positives that are correctly identified by classification technique which is calculated by dividing TP by total of TP and FN as shown in equation (2). F-measure is showing the performance of a classification algorithm and is calculated by the harmonic mean of precision and recall as shown in equation (3).

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

$$F - measure = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (3)$$

We will also report *Receiver Operating Characteristic (ROC)* that is a potentially powerful metric for comparison of different classifiers, because it is invariant against skewness of classes in the dataset. In a *ROC* curve the true positive rate is plotted in function of the false positive rate for different thresholds. In addition to *ROC*, *Area Under the Curve (AUC)* is a measure of how well a parameter can be used to distinguish between two classes. *AUC* is a single value that summarizes the *ROC* by calculating the area of the convex shape below the *ROC* curve. *AUC* can be between 0 and 1, where the value of 1 shows optimal point of perfect prediction.

Matthews Correlation Coefficient (MCC) [17] provides another measures of quality to compare different classifiers [18]. The *MCC* value is between -1 and $+1$, where in cases of perfect prediction it gives $+1$. -1 coefficient shows total disagreement between prediction and observation while the coefficient value of 0 indicates that the classifier does not work better than a random prediction. *MCC* is also a useful measure of classifier performance against imbalanced datasets. While *Precision*, *Recall* or *F-measure* values in a random guessing would be higher than 0.5 , *MCC* value would be around 0 for random guessing. Therefore, for making sure that our classifiers are far from random classifiers, we will compute *MCC* values for each classifier. The values can be computed using equation (7) which is composed of equations (4), (5) and (6), where N is the total number of samples.

$$N = TP + FP + TN + FN \quad (4)$$

$$S = \frac{TP + FN}{N} \quad (5)$$

$$P = \frac{TP + FP}{N} \quad (6)$$

$$MCC = \frac{\frac{TP}{N} - S \times P}{\sqrt{PS(1-S)(1-P)}} \quad (7)$$

The remainder of this paper is organized as follows. Section II reviews some related research in while Section III explains our method for collecting and preprocessing of data in a controlled environment. We describe feature extraction and vectorization in Section IV. Section V introduces our approach for ransomware detection followed by Section VI that describes our performance in detecting ransomwares families. Finally, section VII discusses about the achievements of this paper and concludes the paper.

II. RELATED WORK

Ransomware programs are reportedly becoming a dominant tool for cybercriminals and a growing threat to our ICT infrastructure [4], [19], [20]. The possibility of using encryption techniques to encrypt users data as part of a Denial of Service (DoS) attack is known for a very long time [21]. However, recent adoption of eCurrencies such as BitCoin provided many new opportunities for attackers including receiving a ransom payment for decrypting users data [21]. In spite of its simplicity and primitive utilization of cryptographic techniques [22], ransomware programs are becoming a major tool in cyber criminals toolset [23]. For any cyber threat, prevention is ideal but detection is a must and ransomware is not an exception [3], [24].

Situational cyber security awareness plays an important role in preventing cyber-attacks [25]. An educational framework that is tailored to ransomware threats [10] as well as a tool which mimicked ransomware attacks [26] proved to be useful in reducing ransomware infections. Moreover, technical countermeasures such verifying applications trustworthiness when calling a crypto library [27] or minimizing attack surface by limiting end-users privilege proved effective in preventive ransomware attacks [9].

Most ransomwares detection solutions are relying on filesystem [28]–[30] and registry events [31] to identify malicious behaviors. Investigation of 1359 ransomware samples showed that majority of ransomware samples are using similar APIs and generating similar logs of filesystem activities [29]. For example, using 20 types of filesystem and registry events as features of a *Bayesian Network* model against 20 Windows ransomware samples resulted to an accurate ransomware detection with *F-Measure* of 0.93 [31]. *UNVEIL* [29] as a ransomware classification system utilized filesystem events to distinguish 13,637 ransomwares from a dataset of 148,223 malware samples with accuracy of 96.3%. *CloudRPS* [32] was a cloud-based ransomware detection system which relied on abnormal behaviors such as conversion of large quantities of files in a short interval to detect ransomware samples. *EldeRan* [13] utilized association between different operating system events to build a matrix of applications activities and to detect ransomware samples within 30 seconds of their execution with *AUC* of 0.995. Timely detection of a ransomware upon its execution is very crucial and systems that fail to detect ransomware in less than 10 seconds are not considered effective [5]. Moreover, timely identification of a ransomware family

would assist in building intelligence about applicable threat actors and threat profile for a given target.

III. DATA CREATION

We have downloaded 1624 Windows Portable Executable (PE32) ransomware samples from *virustotal.com* which were active in the period of February 2016 to March 2017 as reported by *RansomwareTracker.abuse.ch*. Collected samples belong to three families of ransomware namely 517 *Locky* samples, 535 *Cerber* samples and 572 samples of *TeslaCrypt*. The best type of goodwill counterpart for malware applications are portable and standalone benign apps [25]. Therefore, we have collected all 220 available portable Windows PE32 benign applications from *portableapps.com*¹ in April 2017 to serve as goodwill counterpart of our dataset.

We have setup the environment shown in Fig. 2 to collect logs of ransomware and goodwill samples runtime activities. The Controller application on the host machine is randomly selecting a ransomware or goodwill sample and passes it through FTP server to the Virtual Machine (VM). When the sample is successfully transferred, the Controller notifies the Launcher app to run the *ProcessMonitor* application and executes a given sample. Similar to the previous research [5], the first 10 seconds log of ransomware and benign applications runtime activities is collected and the created log file is uploaded to the Log repository on the host machine. Since majority of benign applications require human interactions to run (i.e clicking on a button), we have developed an application called *PyWinMonkey* which automates user interactions with an application. When the log file is successfully stored on the host machine, the Controller application reverts the VM back to its original copy and passes the next sample. It is notable that *PyWinMonkey* is similar to *Monkey*² Android app which utilized in many previous Android malware research papers [33] for mimicking human interactions. We have used Python 3.6.1 to develop Controller, Launcher and *PyWinMonkey* apps (accessible at <https://github.com/sajadhomayoun/PyWinMonkey>) and run *ProcessMonitor V3.31* on Windows10 build number 10240 on a computer with Core i7 CPU with 8 cores of 4GHz and 16GB of RAM. For each and every process, *ProcessMonitor* records loaded Dynamic Linked Libraries (DLLs), file system activities and registry activities. Therefore, we will have three sets of events namely *Registry_Events_Set*, which includes all registry events, *DLL_Events_Set*, which includes all DLL events and *Filesystem_Events_Set*, which contains all Filesystem events as listed in Table I. Moreover, *EventType(E)* is a procedure that returns the type of given event (R for Registry events, F for Filesystem events, and D for DLL events) as Fig. 1.

As we will be using a sequential pattern mining technique (MG-FSM) to detect candidate features for classification task, we should convert our data into a sequential dataset which is a collection of sequences such as $D = \{S_1, S_2, \dots, S_n\}$ where S_i represents a sequentially ordered set of events. We have created a sequence of runtime events for each and every

```

1: procedure EventType(Event E)
2:   if  $E \in Registry\_Events\_Set$  return R
3:   if  $E \in Filesystem\_Events\_Set$  return F
4:   if  $E \in DLL\_Events\_Set$  return D
5: end procedure

```

Fig. 1. Determining Even Type of a given event.

TABLE I
LIST OF ACTIVITIES CAN BE CAPTURED BY PROCESS MONITOR

Activity Type	List
Registry	RegQueryKey, RegOpenKey, RegQueryValue, RegCloseKey, RegCreateKey, RegSetInfoKey, RegEnumKey, RegQueryKeySecurity, RegEnumValue, RegSetValue, RegDeleteValue, RegQueryMultipleValueKey, RegDeleteKey, RegLoadKey, RegFlushKey
File	QueryNameInformationFile, ReadFile, CreateFile, QueryBasicInformationFile, CloseFile, QueryStandardInformationFile, CreateFileMapping, QuerySizeInformationVolume, FileSystemControl, QueryDirectory, WriteFile, QueryNetworkOpenInformationFile, QueryRemoteProtocolInformation, QuerySecurityFile, LockFile, UnlockFileSingle, DeviceIoControl, SetEndOfFileInformationFile, FlushBuffersFile, SetAllocationInformationFile, SetBasicInformationFile, QueryAttributeTagFile, QueryFileInternalInformationFile, QueryInformationVolume, QueryAttributeInformationVolume, SetRenameInformationFile, QueryNormalized-NameInformationFile, NotifyChangeDirectory, QueryFullSizeInformationVolume, SetSecurityFile, QueryStreamInformationFile, SetDispositionInformationFile, QueryEaInformationFile, QueryAllInformationFile, QueryIdInformation, SetPositionInformationFile, QueryPositionInformationFile, SetValidDataLengthInformationFile
DLL	LoadImage

ransomware and benign application. S_i represents a sequence of all events E caused by launching an application i ordered by time as follow:

$S_i = \{E_{1,i}(argE_1), E_{2,i}(argE_2), \dots, E_{n,i}(argE_n)\}$ where $E_{x,y}(argE_x)$ represents event x for an application y and $argE_x$ shows the argument passed to the event E_x .

For example, $\{LoadImage(C : \backslash system32 \backslash gdi32.dll)\}$, $\{LoadImage(ReadFile(C : \backslash Windows \backslash SysWOW64 \backslash wininet.dll))\}$ shows a sequence of two events where the first event loads *gdi32.dll* in the memory of calling process (hence $C : \backslash system32 \backslash gdi32.dll$ is the parameter for this event) and the second event reads *wininet.dll* file located at $C : \backslash Windows \backslash SysWOW64$. The size of each sequence depends on the number of events that are called by an application and varies between different apps.

Once all sequences are created, we have utilized the *Outlier Factor* [34] technique to remove any outlier sequence from our dataset similar to [35]. The *Outlier Factor* technique first extracts all frequent patterns from a dataset and then detects

¹<https://portableapps.com/apps>

²<https://developer.android.com/studio/test/monkey.html>

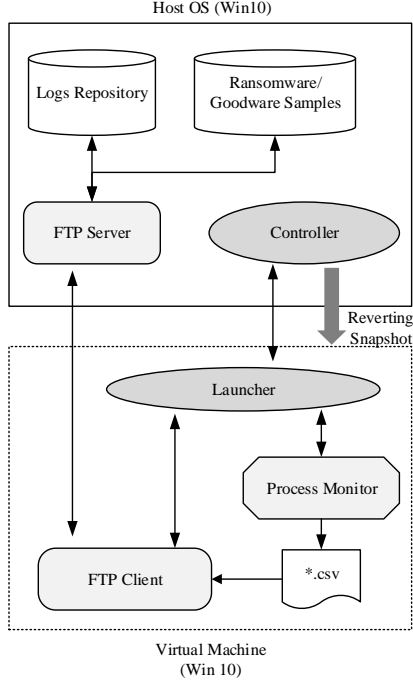


Fig. 2. Environment Setup to Capture Malware and Goodware Activities Log

TABLE II
CREATED DATASETS

Dataset	Number of Sequences
D_Locky	450
D_Cerber	470
D_TeslaCrypt	507
D_Goodware	200
D_OF	174

outlier sequences as those that contain the least frequent patterns in a given dataset.

Table II reflects final datasets with the number of sequences in each dataset. *D_Locky* represents sequences of *Locky* ransomware samples, *D_Cerber* shows *Cerber* ransomware sequences and *D_TeslaCrypt* includes sequences of *TeslaCrypt* ransomware samples. *D_Ransomware* represents combined sequences of all ransomware samples while *D_Goodware* includes sequences of events of all benign applications. We randomly collected 52 *Locky*, 50 *Cerber*, 52 *TeslaCrypt* and 20 benign applications sequences in a separated dataset for over-fitting test as well (*D_OF*).

IV. FEATURE EXTRACTION AND VECTORIZATION

To detect the best features for classification task, we need to first define detectable patterns of events and then utilize a pattern mining algorithm to find *Maximal Sequential Patterns (MSP)* collections within each dataset. Afterwards, every sequence within every relevant dataset is traversed based on a given *MSP* collection to provide features for training classifiers.

Sequential pattern mining techniques discover all subsequences (Sequential Patterns) that appear in a given sequential dataset with frequency of no less than a user-specified threshold (min_{sup}) [36]. A sequence $\alpha = \{a_1, a_2, \dots, a_n\}$ is called a subsequence of another sequence $\beta = \{b_1, b_2, \dots, b_m\}$ and β is a super-sequence of α , denoted as $\alpha \subseteq \beta$, if there exists integers $1 \leq j_1 < j_2 < \dots < j_n \leq m$ such that $a_1 \subseteq b_{j_1}, a_2 \subseteq b_{j_2}, \dots, a_n \subseteq b_{j_n}$. A sequence is said to be frequent and called a *Sequential Pattern (SP)* in a sequential dataset D if $sup_{\alpha} \geq min_{sup}$, where sup_{α} (support of α) denotes the frequency of occurrence of α in a given sequential dataset D . Moreover, if a *Sequential Pattern SP* is not contained in any other sequential patterns, it is called a *Maximal Sequential Pattern (MSP)*. Collection of all *MSPs* within a given sequential dataset D can be denoted as a *Maximal Sequential Pattern Collection (MC_D)*. Members of a *MC* are in format of (P, sup_P) where P is a *MSP* and sup_P shows the frequency of occurrence of P in a given dataset D .

There are two major types of sequential pattern mining algorithms to extract *MSPs* namely Apriori-based and frequent pattern growth. Apriori-based algorithms are detecting *MSPs* based on the fact that any subset of a frequent pattern must be frequent. However, recursive nature of Apriori-based algorithms increases complexity and running time of the algorithm [37]. On the other side, frequent pattern growth algorithms are using divide-and-conquer techniques to narrow down the search space *MSPs*. To detect *MSPs* in this study, we utilize a widely used frequent pattern growth algorithm [38] called "Mind the Gap: Frequent Sequence Mining (MG-FSM)" [39] with min_{sup} of 50%. Applying MG-FSM against our datasets generates four *MSP* collections namely MC_{D_Locky} , MC_{D_Cerber} , $MC_{D_TeslaCrypt}$ and $MC_{D_Ransomware}$.

$$MC_D = \{(P_x, sup_{P_x}) | sup_{P_x} \geq min_{sup} \wedge \forall P_y (P_x \subseteq P_y) \}.$$

We can distinguish three types of atomic *MSPs* and six types of single step transition *MSPs* within our sequential datasets as shown in Table III. Atomic *MSPs* are representing continuous events of the same type i.e. the atomic *MSP* of F represents continuous Filesystem events. Single step transitions *MSPs* are representing a transition from one atomic *MSP* to another. For example, *MSP* of RD represents a sequence of registry events (R atomic *MSP*) followed by a sequence of DLL events (D atomic *MSP*). It is notable that we only define two types of atomic and single step transition *MSP* to avoid sparsity in extracted features.

A *MSP* $P = \{E_1, \dots, E_n\}$ is atomic if $\forall E_x, E_y \in P \wedge E_x \neq E_y (EventType(E_x) == EventType(E_y))$.

A *MSP* $P = \{E_1, \dots, E_n\}$ is a single step transition if $\exists E_x, E_y \in P \wedge E_x \neq E_y (EventType(E_x) \neq EventType(E_y))$.

We can define a set that contains all *MSP* types (*MSP Type_Set*) and a procedure (*MSPTYPE(MSP P)*) in Fig. 3) that returns type of given sequence S as follow:

$$MSP_Type_Set = \{R, F, D, RF, RD, FR, FD, DR, DF\}.$$

Support Ratio (SR) of a *MSP* is a value in the range of [0,1] that shows the possibility of occurrence of the *MSP* in a given dataset of ransomware and is calculated by dividing

```

1: procedure MSPTYPE(MSP P)
2:   for all  $(E_x, E_y \in P) \wedge (x \leq i) \wedge (y > i) \wedge (i, y \leq n)$ 
3:   do
4:     if  $EventType(E_x) == EventType(E_y)$  then
5:       if  $EventType(E_x) == R$  return R
6:       if  $EventType(E_x) == F$  return F
7:       if  $EventType(E_x) == D$  return D
8:     else
9:       if  $EventType(E_x) == R \wedge$ 
10:       $EventType(E_y) == F$  return RF
11:      if  $EventType(E_x) == R \wedge$ 
12:       $EventType(E_y) == D$  return RD
13:      if  $EventType(E_x) == F \wedge$ 
14:       $EventType(E_y) == R$  return FR
15:      if  $EventType(E_x) == F \wedge$ 
16:       $EventType(E_y) == D$  return FD
17:      if  $EventType(E_x) == D \wedge$ 
18:       $EventType(E_y) == R$  return DR
19:      if  $EventType(E_x) == D \wedge$ 
20:       $EventType(E_y) == F$  return DF
21:     end if
22:   end for
23: end procedure

```

Fig. 3. Finding MSP Type of a given MSP.

TABLE III
MAXIMAL SEQUENTIAL PATTERN TYPES

Type	Description
R	All events must be registry
F	All events must be file
D	All events must be actions on dll files
RF	The MFP has one or more transitions while the first transition is from a registry event to a file event
RD	The MFP has one or more transitions while the first transition is from a registry event to a dll event
FR	The MFP has one or more transitions while the first transition is from a file event to a registry event
FD	The MFP has one or more transitions while the first transition is from a file event to a dll event
DR	The MFP has one or more transitions while the first transition is from a dll event to a registry event
DF	The MFP has one or more transitions while the first transition is from a dll event to a file event

frequency of occurrences of *MSP* (sup_MSP) by the total number of all ransomware sequences (γ) in a given dataset D . For every sequence S we can define a Vector of size nine (9) that contains *SR* value of every *MSP* type detected in *MSP Collection* MC within sequence S as follow:

$$Vector(S)_{MC} = \{(SR_R), (SR_F), (SR_D), (SR_RF), (SR_RD), (SR_FR), (SR_FD), (SR_DR), (SR_DF)\}$$

SR value of every *MSP Type* of a sequence for a given MC can be calculated using *CalculateSR* procedure shown in Fig 4. When vector of all sequences within a sequential dataset

```

1: procedure CALCULATESR(Sequence S, MSP Collection
2:   MC, MSP_Type_Set T)
3:    $SR\_P\_Total = 0$ 
4:   for all  $P \in MC$  do
5:     if  $p \subseteq S$  AND  $MSPTYPE(P) == T$  then
6:        $SR\_P\_Total = SR\_P\_Total + (\frac{sup\_P}{\gamma})$ 
7:     end if
8:   end for
9:   return  $SR\_P\_Total$ 
10: end procedure

```

Fig. 4. SR calculation algorithm for sequence S.

D using a *MSP Collection* MC is created, we will have a Vektored Dataset $VD_{D,MC}$.

Moreover, for every sequence S we can define a *SuperVector* of S as a set of Vectors created using *MSPs* collected from different dataset i.e. $MC_{D_1}, MC_{D_2}, \dots, MC_{D_n}$ as follow:

$$SuperVector(S) = \{Vector(S)_{MC_1}, Vector(S)_{MC_2}, \dots, Vector(S)_{MC_n}\}$$

where MC_1 to MC_n reflects *MSP Collections* of different datasets. Size of a *SuperVector* with n vectors is calculated by $n \times m$ where m is size of each vector (9 in this research). By calculating *SuperVector* for all sequences within a dataset D , we will have Super Vektored Dataset $SV_{D,D}$.

V. HUNTING FOR EVIL: RANSOMWARE DETECTION

To detect best features (*MSP* types) for classifying ransomware from benign applications we created a dataset MC_{D_Total} by combining $D_Ransomware$ and $D_Goodware$ and then generated a vektored dataset $VD_{D_Total,MC_{Ransomware}}$. We then utilized greedy stepwise search method of *CfsSubsetEval* [40] of *Weka3.8.1* with $VD_{D_Total,MC_{Ransomware}}$ and found that *MSP Types* of R (Registry), D (DLLs) and FD (File and DLL) may provide best distinction between ransomware and goodware samples (see Fig. 5). As shown in Fig. 5a ransomware applications are tend to conduct a much wider range of Registry activities in compare with goodware apps. As shown in Fig. 5b, majority of benign applications were conducting similar DLL activities while there were much more variations in ransomware samples DLL events. Ransomware applications are taking a variety of Filesystem to DLL transitions while goodware samples were mainly taking only two specific Filesystem to DLL events transitions (see Fig. 5c).

We have utilized $R, D,$ and FD as features to train four classifiers namely *J48, Random Forest, Bagging,* and *Multi Layer Perception (MLP)* using $VD_{D_Total,MC_{Ransomware}}$ and 10-fold cross validation technique for evaluation. As shown in Table IV, all classifiers achieved F-measure of 0.99 with a low false positive rate ($FPR \leq 0.04$). Moreover, similarities between *ROC* curves of different classifiers (see Fig 6) proves that there is not much difference between performance of different classifiers which is another indication of suitability of our features for classifying ransomware and benign applications. As shown in Fig 7, *AUC* value for all classifiers is quite high (more than 0.990) while *AUC* value of *Bagging* classifier

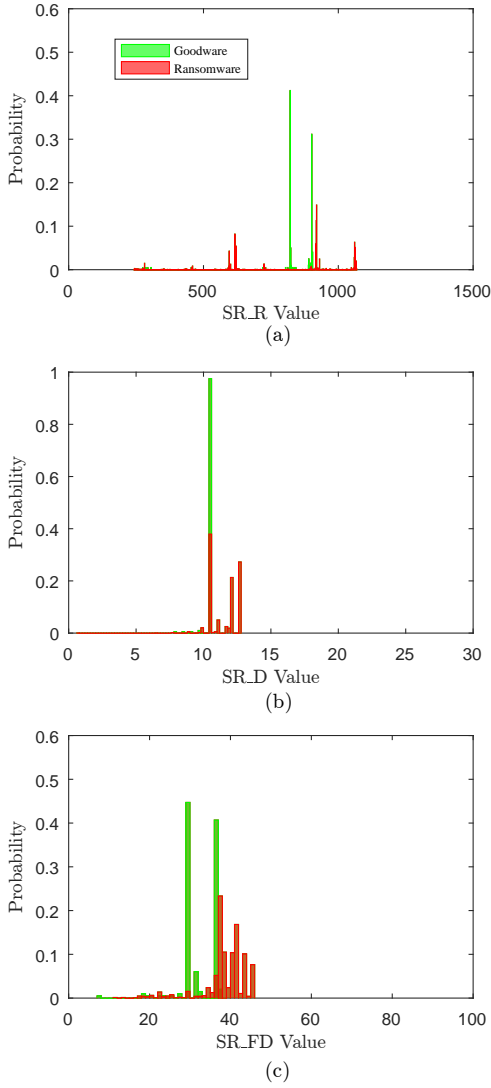


Fig. 5. Histogram of the probability of SR values for ransomware and goodware.

TABLE IV
CLASSIFIERS PERFORMANCE ON $V D_{D_Total, MC_{Ransomware}}$

Classifier	TPR	FPR	F-Measure
J48	0.994	0.040	0.994
Random Forest	0.993	0.040	0.993
Bagging	0.994	0.039	0.977
MLP	0.994	0.035	0.994

(0.995) is very close to an optimal prediction. The *MCC* value of all classifiers is more than 0.96 while *Random Forest* and *Bagging* achieved *MCC* of almost +1 which is very close to a perfect prediction.

To show that we have not over-fitted our classifiers, we tested all classifiers using on $V D_{D_OF, MC_{Ransomware}}$. As

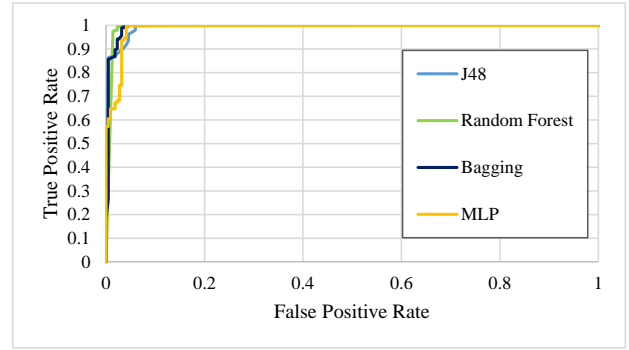


Fig. 6. ROC diagrams for classifiers.

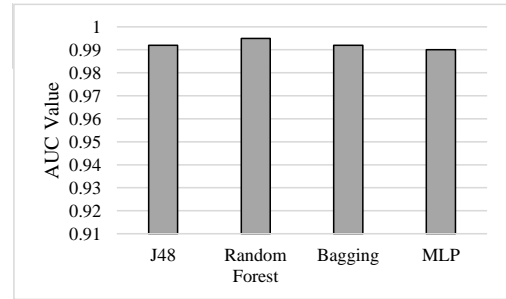


Fig. 7. AUC of classifiers for detecting ransomwares

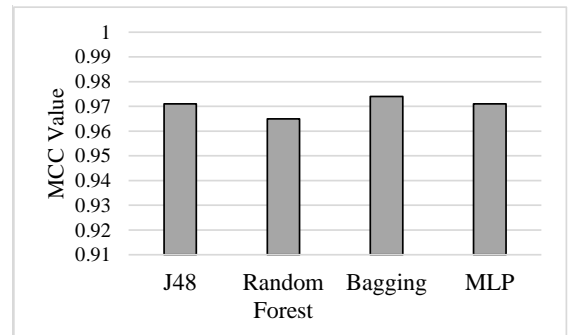


Fig. 8. MCC of classifiers for detecting ransomwares

shown in Table V all classifiers achieved accuracy of 0.994 in classifying unforeseen ransomware and goodware samples too.

TABLE V
RESULTS OF CLASSIFIERS ON $VD_{D_OF,MC_{Ransomware}}$ FOR DETECTING RANSOMWARE

Classifier	Accuracy
J48	0.994
Random Forest	0.994
Bagging	0.994
MLP	0.994

VI. THREAT INTELLIGENCE: DETECTION OF A RANSOMWARE FAMILY

To investigate performance of classifiers in detection of a ransomware family we have created $D_{TotalFamily}$ dataset which contains all sequences from D_{Locky} , D_{Cerber} , $D_{TeslaCrypt}$ and $D_{Goodware}$. We then generated $VD_{D_{TotalFamily},MC_{Locky}}$, $VD_{D_{TotalFamily},MC_{Cerber}}$ and $VD_{D_{TotalFamily},MC_{TeslaCrypt}}$ vectored datasets and fed them to *CfsSubsetEval* of *Weka3.8.1*. All in all 13 candidate features were detected for classification of ransomware families as shown in Fig. 9.

Fig. 9a depicts the distribution of SR_R , SR_{RF} , SR_{FD} and SR_{DF} values in $VD_{D_{TotalFamily},MC_{Locky}}$. The interesting point in Fig. 9a is that *Locky* samples have taken higher values for all selected features in $VD_{D_{TotalFamily},MC_{Locky}}$ in compare with *Cerber*, *TeslaCrypt* and *goodware* samples. This is because $VD_{D_{TotalFamily},MC_{Locky}}$ is vectored using MC_{Locky} , so it is vivid that more *MSPs* are matched to *Locky* samples which is reflected in higher *SR* values for corresponding features. Fig. 9a also shows that behaviours of *Cerber* and *TeslaCrypt* samples have been different from *Locky*, and this leads to smaller values for features R, RF, FD and DF.

Fig. 9b shows histograms of distribution of *SR* values namely SR_R , SR_D , SR_{FR} , SR_{FD} and SR_{DF} in $VD_{D_{TotalFamily},MC_{Cerber}}$. The histograms in Fig. 9b show higher *SR* values for *Cerber* samples in compare with samples of other classes and it means *Cerber* samples matched more *MSPs* from MC_{Cerber} .

Fig. 9c illustrates the histogram of *SR* values for SR_R , SR_F , SR_{RF} and SR_{FD} for $VD_{D_{TotalFamily},MC_{TeslaCrypt}}$. The diagrams in Fig. 9c present higher *SR* values for *TeslaCrypt* samples in compare with *Locky*, *Cerber* and *goodware* samples. Fig. 9 also demonstrates that *goodware* samples always take smaller and far *SR* values for each feature because *goodwares* behave differently in compare with *ransomwares*. In other words, *goodware* samples have matched fewer *MSPs* from MC_{Locky} , MC_{Cerber} and $MC_{TeslaCrypt}$, and subsequently they have taken lower *SR* values.

As detection of ransomware families is a multi-class classification task with four class labels (*Locky*, *Cerber*, *TeslaCrypt* and *Goodware*), therefore, we have trained J48, *Random Forest*, *Bagging* and *MLP* with a multi-class classifier using $SV_{D_{D_{TotalFamily}}}$ dataset with 13 selected features in Fig. 9.

Table VI presents performance of all classifiers obtained from 10-fold cross validation. Obtained minimum weighted

TABLE VI
THE CLASSIFIERS PERFORMANCE ON $SV_{D_{D_{TotalFamily}}}$

Classifier	TPR	FPR	F-Measure	MCC
J48	0.981	0.006	0.981	0.974
Random Forest	0.983	0.006	0.983	0.978
Bagging	0.980	0.007	0.980	0.974
MLP	0.980	0.007	0.980	0.973

TABLE VII
RESULTS OF CLASSIFIERS ON DATASET $SV_{D_{D_OF}}$ FOR DETECTING RANSOMWARE FAMILY

Classifier	Accuracy
J48	0.947
Random Forest	0.965
Bagging	0.959
MLP	0.959

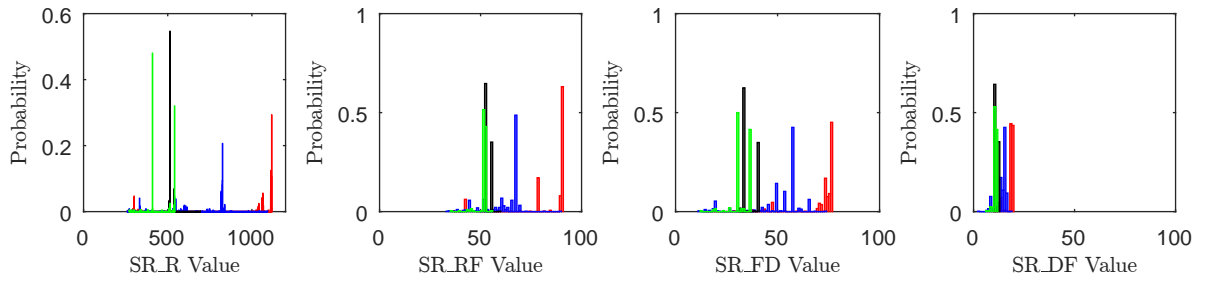
average [41] *F-Measure* of 0.983 with $FPR \leq 0.006$ reflects suitability of our features for detecting ransomware samples families. *MCC* values of more than 0.95 for all classifiers also indicate quality of our features in enabling classifiers to provide an almost perfect prediction. Finally, as shown in Table VII our features enabled classifiers to offer an accurate prediction (≥ 0.965) even on unforeseen samples ($SV_{D_{D_OF}}$).

VII. CONCLUDING REMARKS

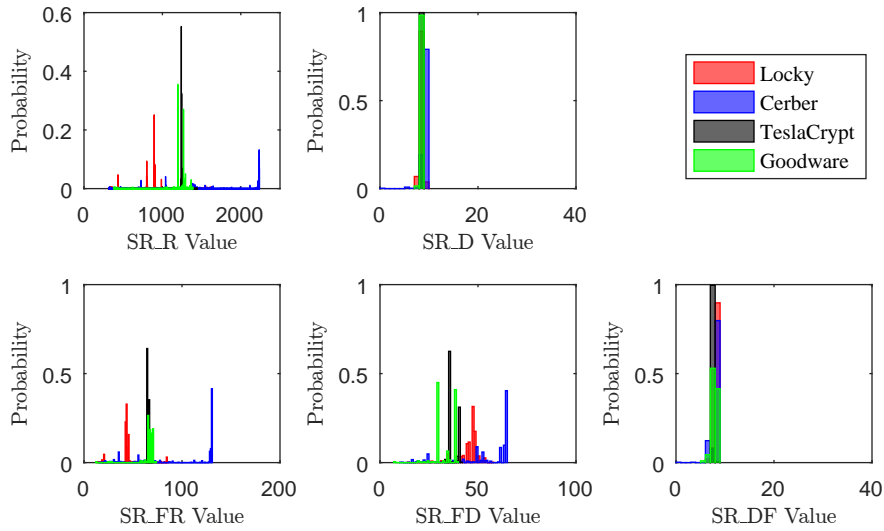
In this paper, by combining sequential pattern mining for feature identification with machine learning classification techniques we could accurately distinguish between ransomware and *goodware* samples and identify given ransomware families with in first 10 seconds of a ransomware execution. We achieved minimum *F-measure* of 0.994 with minimum *AUC* value of 0.99 in detection of ransomware samples from *goodware* using Registry (R) events, DLL (D) events and Filesystem to Registry (FD) transitions as features for J48, *Random Forest*, *Bagging* and *MLP* classifiers. We achieved *F-Measure* of more than 0.98 with *FPR* of less than 0.007 in detection of a given ransomware family using 13 selected features detected in this study. Reported features for differentiating ransomware and benign applications can be used for effective hunting of a ransomware while features reported for ransomware family classification are great for building intelligence about threat profiles applicable to a given target. Applying other classification techniques such as fuzzy classification can be considered as a future work of this study. Moreover, utilization of *Stream Data Mining* techniques to reduce ransomware detection time is another interesting extension of this study.

ACKNOWLEDGMENT

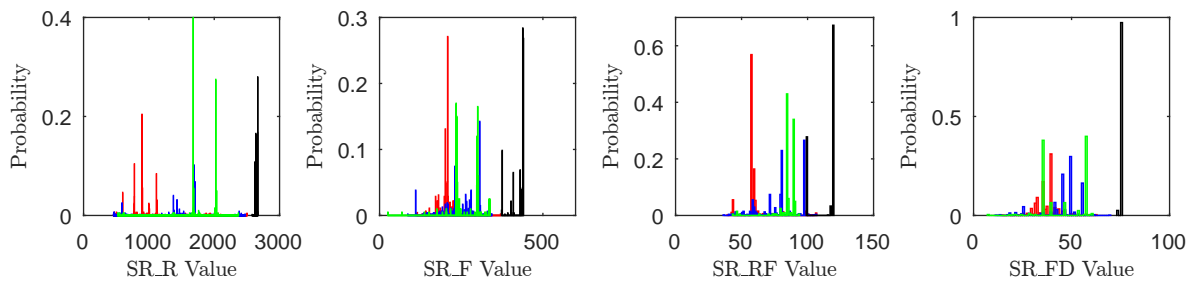
The authors would like to thank *virustotal.com* for sharing malware samples and giving the capability of scanning files. We also thank *ransomwaretracker.abuse.ch* for their updated and precious information about different families of



(a) Features selected from $VD_{D_TotalFamily,MC_{Locky}}$



(b) Features selected from $VD_{D_TotalFamily,MC_{Cerber}}$



(c) Features selected from $VD_{D_TotalFamily,MC_{TeslaCrypt}}$

Fig. 9. Histogram of the probability of SR values for ransomware families.

ransomwares. This work is partially supported by the European Council 268 International Incoming Fellowship (FP7-PEOPLE-2013-IIF) grant.

REFERENCES

- [1] M. Hopkins and A. Dehghantanha, "Exploit kits: The production line of the cybercrime economy?" in *2015 Second International Conference on Information Security and Cyber Forensics (InfoSec)*. IEEE, nov 2015. [Online]. Available: <https://doi.org/10.1109%2Finfosec.2015.7435501>
- [2] EUROPOL. (2016) The internet organised crime threat assessment (iocta) 2016. [Online]. Available: <https://www.europol.europa.eu/activities-services/main-reports/internet-organised-crime-threat-assessment-iocta-2016>
- [3] N. Milosevic, A. Dehghantanha, and K.-K. R. Choo, "Machine learning aided android malware classification," *Computers & Electrical Engineering*, feb 2017. [Online]. Available: <https://doi.org/10.1016%2Fj.compeleceng.2017.02.013>
- [4] K. Cabaj and W. Mazurczyk, "Using software-defined networking for ransomware mitigation: The case of CryptoWall," *IEEE Network*, vol. 30, no. 6, pp. 14–20, nov 2016. [Online]. Available: <https://doi.org/10.1109/nnet.2016.1600110nm>
- [5] A. Azmoodeh, A. Dehghantanha, M. Conti, and R. Choo, "Detecting crypto-ransomware in iot networks based on energy consumption footprint," *Journal of Ambient Intelligence and Humanized Computing*, 2017.
- [6] H. L. Kevin Savage, Peter Coogan, *The evolution of ransomware*. Symantec, 2015.
- [7] Symantec, "Internet security threat report," Symantec, Tech. Rep., apr 2016.
- [8] K. T. D. Y. Huang, D. W. E. B. C. GrierD, T. J. Holt, C. Kruegel, D. McCoy, S. Savage, and G. Vigna, "Framing dependencies introduced by underground commoditization," in *Workshop on the Economics of Information Security*, 2015.
- [9] Monika, P. Zavorsky, and D. Lindskog, "Experimental analysis of ransomware on windows and android platforms: Evolution and characterization," *Procedia Computer Science*, vol. 94, pp. 465–472, 2016. [Online]. Available: <https://doi.org/10.1016/j.procs.2016.08.072>
- [10] X. Luo and Q. Liao, "Awareness education as the key to ransomware prevention," *Information Systems Security*, vol. 16, no. 4, pp. 195–202, sep 2007. [Online]. Available: <https://doi.org/10.1080%2F10658980701576412>
- [11] M. Simmonds, "How businesses can navigate the growing tide of ransomware attacks," *Computer Fraud & Security*, vol. 2017, no. 3, pp. 9–12, mar 2017. [Online]. Available: [https://doi.org/10.1016/s1361-3723\(17\)30023-4](https://doi.org/10.1016/s1361-3723(17)30023-4)
- [12] S. Grzonkowski, A. Mosquera, L. Aouad, and D. Morss, "Smartphone security: An overview of emerging threats," *IEEE Consumer Electronics Magazine*, vol. 3, no. 4, pp. 40–44, oct 2014. [Online]. Available: <https://doi.org/10.1109/mce.2014.2340211>
- [13] D. Sgandurra, L. Muoz-Gonzlez, R. Mohsen, and E. C. Lupu, "Automated dynamic analysis of ransomware: Benefits, limitations and use for detection," 2016.
- [14] M. Sohrabi, M. M. Javidi, and S. Hashemi, "Detecting intrusion transactions in database systems: a novel approach," *Journal of Intelligent Information Systems*, vol. 42, no. 3, pp. 619–644, dec 2013. [Online]. Available: <https://doi.org/10.1007%2Fs10844-013-0286-z>
- [15] M. Sun, X. Li, J. C. S. Lui, R. T. B. Ma, and Z. Liang, "Monet: A user-oriented behavior-based malware variants detection system for android," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 5, pp. 1103–1112, may 2017. [Online]. Available: <https://doi.org/10.1109%2Ftifs.2016.2646641>
- [16] M. R. Watson, N. ul-hassan Shirazi, A. K. Marnerides, A. Mauthe, and D. Hutchison, "Malware detection in cloud computing infrastructures," *IEEE Transactions on Dependable and Secure Computing*, vol. 13, no. 2, pp. 192–205, mar 2016. [Online]. Available: <https://doi.org/10.1109%2Ftdsc.2015.2457918>
- [17] B. Matthews, "Comparison of the predicted and observed secondary structure of t4 phage lysozyme," *Biochimica et Biophysica Acta (BBA) - Protein Structure*, vol. 405, no. 2, pp. 442–451, oct 1975. [Online]. Available: <https://doi.org/10.1016%2F0005-2795%2875%2990109-9>
- [18] D. M. Powers, "Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation," School of Informatics and Engineering-Flinders University, Tech. Rep., 2011.
- [19] "Ransomware becomes most popular form of attack as payouts approach \$1bn a year," *Network Security*, vol. 2017, no. 1, pp. 1–2, jan 2017. [Online]. Available: [https://doi.org/10.1016/s1353-4858\(17\)30001-6](https://doi.org/10.1016/s1353-4858(17)30001-6)
- [20] "UK major target for ransomware," *Computer Fraud & Security*, vol. 2016, no. 1, p. 3, jan 2016. [Online]. Available: [https://doi.org/10.1016/s1361-3723\(16\)30003-3](https://doi.org/10.1016/s1361-3723(16)30003-3)
- [21] H. Orman, "Evil offspring - ransomware and crypto technology," *IEEE Internet Computing*, vol. 20, no. 5, pp. 89–94, sep 2016. [Online]. Available: <https://doi.org/10.1109/mic.2016.90>
- [22] C. Everett, "Ransomware: to pay or not to pay?" *Computer Fraud & Security*, vol. 2016, no. 4, pp. 8–12, apr 2016. [Online]. Available: [https://doi.org/10.1016/s1361-3723\(16\)30036-7](https://doi.org/10.1016/s1361-3723(16)30036-7)
- [23] A. Gazet, "Comparative analysis of various ransomware virii," *Journal in Computer Virology*, vol. 6, no. 1, pp. 77–90, jul 2008. [Online]. Available: <https://doi.org/10.1007%2Fs11416-008-0092-2>
- [24] R. Brewer, "Ransomware attacks: detection, prevention and cure," *Network Security*, vol. 2016, no. 9, pp. 5–9, sep 2016. [Online]. Available: [https://doi.org/10.1016/s1353-4858\(16\)30086-1](https://doi.org/10.1016/s1353-4858(16)30086-1)
- [25] M. Damshenas, A. Dehghantanha, and R. Mahmoud, "A survey on malware propagation, analysis, and detection," *International Journal of Cyber-Security and Digital Forensics (IJCSDF)*, vol. 2, no. 4, pp. 10–29, 2013.
- [26] S. Al-Sharif, F. Iqbal, T. Baker, and A. Khattack, "White-hat hacking framework for promoting security awareness," in *2016 8th IFIP International Conference on New Technologies, Mobility and Security (NTMS)*. IEEE, nov 2016. [Online]. Available: <https://doi.org/10.1109/ntms.2016.7792489>
- [27] A. L. Young, "Cryptoviral extortion using microsoft crypto API," *International Journal of Information Security*, vol. 5, no. 2, pp. 67–76, mar 2006. [Online]. Available: <https://doi.org/10.1007%2Fs10207-006-0082-7>
- [28] A. Kharraz, W. Robertson, D. Balzarotti, L. Bilge, and E. Kirida, "Cutting the gordian knot: A look under the hood of ransomware attacks," in *Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer Nature, 2015, pp. 3–24. [Online]. Available: https://doi.org/10.1007%2F978-3-319-20550-2_1
- [29] A. Kharaz, S. Arshad, C. Mulliner, W. Robertson, and E. Kirida, "Unveil: A large-scale, automated approach to detecting ransomware," in *25th USENIX Security Symposium (USENIX Security 16)*. Austin, TX: USENIX Association, 2016, pp. 757–772. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/kharaz>
- [30] C. Moore, "Detecting ransomware with honeypot techniques," in *2016 Cybersecurity and Cyberforensics Conference (CCC)*. Institute of Electrical and Electronics Engineers (IEEE), aug 2016. [Online]. Available: <https://doi.org/10.1109%2Fccc.2016.14>
- [31] M. M. Ahmadian and H. R. Shahriari, "2entfox: A framework for high survivable ransomwares detection," in *2016 13th International Iranian Society of Cryptology Conference on Information Security and Cryptology (ISCISC)*. Institute of Electrical and Electronics Engineers (IEEE), sep 2016. [Online]. Available: <https://doi.org/10.1109%2Fiscisc.2016.7736455>
- [32] J. K. Lee, S. Y. Moon, and J. H. Park, "CloudRPS: a cloud analysis based enhanced ransomware prevention system," *The Journal of Supercomputing*, jul 2016. [Online]. Available: <https://doi.org/10.1007/s11227-016-1825-5>
- [33] M. Damshenas, A. Dehghantanha, K.-K. R. Choo, and R. Mahmud, "M0droid: An android behavioral-based malware detection model," *Journal of Information Privacy and Security*, vol. 11, no. 3, pp. 141–157, jul 2015. [Online]. Available: <https://doi.org/10.1080/15536548.2015.1073510>
- [34] Z. He, X. Xu, Z. Huang, and S. Deng, "FP-outlier: Frequent pattern based outlier detection," *Computer Science and Information Systems*, vol. 2, no. 1, pp. 103–118, 2005. [Online]. Available: <https://doi.org/10.2298%2Fcsis0501103h>
- [35] J. Yu, G. X. Yao, and W. W. Zhang, "Intrusion detection method based on frequent pattern," *Advanced Materials Research*, vol. 204–210, pp. 1751–1754, feb 2011. [Online]. Available: <https://doi.org/10.4028/www.scientific.net/amr.204-210.1751>
- [36] R. Agrawal and R. Srikant, "Mining sequential patterns," in *Proceedings of the Eleventh International Conference on Data Engineering*, ser. ICDE '95. Washington, DC, USA: IEEE Computer Society, 1995, pp. 3–14. [Online]. Available: <http://dl.acm.org/citation.cfm?id=645480.655281>
- [37] W. Shen, J. Wang, and J. Han, "Sequential pattern mining," in *Frequent Pattern Mining*. Springer International Publishing, 2014, pp. 261–282. [Online]. Available: https://doi.org/10.1007/978-3-319-07821-2_11

- [38] K. Damevski, D. C. Shepherd, J. Schneider, and L. Pollock, "Mining sequences of developer interactions in visual studio for usage smells," *IEEE Transactions on Software Engineering*, vol. 43, no. 4, pp. 359–371, apr 2017. [Online]. Available: <https://doi.org/10.1109/2Ftse.2016.2592905>
- [39] I. Miliaraki, K. Berberich, R. Gemulla, and S. Zoupanos, "Mind the gap," in *Proceedings of the 2013 international conference on Management of data - SIGMOD 13*. Association for Computing Machinery (ACM), 2013. [Online]. Available: <https://doi.org/10.1145/2F2463676.2465285>
- [40] M. A. Hall, "Correlation-based feature selection for machine learning," Tech. Rep., 1998.
- [41] M. Sokolova and G. Lapalme, "A systematic analysis of performance measures for classification tasks," *Information Processing & Management*, vol. 45, no. 4, pp. 427–437, jul 2009. [Online]. Available: <https://doi.org/10.1016/j.ipm.2009.03.002>



Sajad Homayoun is a Ph.D candidate of Computer Networks at Shiraz University of Technology since 2013. He also has a master's degree in Information Technology from K. N. Toosi University of Technology (Khajeh Nasire Toosi) of Tehran and a bachelor's degree in Computer Software. He is currently in charge of security laboratory (SecLab) in Shiraz University of Technology since 2014. His research interests are Cyber Security, Machine Learning Applications in Computer Security and Computer Networks.



Ali Dehghantanha is a Marie-Curie International Incoming Fellow in Cyber Forensics, a fellow of the UK Higher Education Academy (HEA) and an IEEE Sr. member. He has served for many years in a variety of research and industrial positions. Other than Ph.D in Cyber Security he holds several professional certificates such as GXPN, GREM, GCFA, CISM, and CISSP.



Marzieh Ahmadzadeh holds a Ph.D in Computer Science and MSc. In Information Technology, both received from the University of Nottingham, UK, and a first class BSc in Software Engineering received from Isfahan University, Iran. Since September 2006, she has been an assistant professor at the school of computer Engineering and IT, Shiraz University of Technology, where she has supervised more than 20 MSc students whose research area are mostly applied data mining. Being a research assistant at the University of Nottingham and full

stack software engineer for 3 years is also part of her work experience. Her research interest includes Data Mining, Data Security, HCI and Computer Science Education.



Sattar Hashemi received the PhD degree in computer science from Iran University of Science and Technology in conjunction with Monash University, Australia, in 2008. Following academic appointments at Shiraz University, he is currently an associate professor at Electrical and Computer Engineering School, Shiraz University, Shiraz, Iran. His research interests include machine learning, data mining, social networks, data stream mining, game theory, and adversarial learning.



Raouf Khayami received his BS degree in computer engineering (hardware systems) from Shiraz University in 1993, the MS degree in artificial intelligence and robotics from the same university in 1996, and the Ph.D degree in software systems from Shiraz University in 2009. He is currently an assistant professor in the Computer Engineering and Information Technology Department, Shiraz University of Technology, Shiraz, Iran, and there, he is the Head of the Department. His research interests include data mining, business intelligence, and enterprise

architecture, on which he has published a number of refereed articles, surveys and technical reports in prestigious national and international conferences and journals. He also is active in consulting and industrial projects.