Paper: Utilizing Cyber Threat Hunting Techniques to Find Ransomware Attacks: A Survey of the State of the Art

Fatimah Aldauiji, Omar Batarfi, Manal Bayousef

Presented by: Oscar Llerena





Abstract

- Ransomware detection and response needs
- Cyber Threat Hunting (CTH) approaches
- Cyber Threat Intelligence (CTI) techniques
- Paper investigates practical CTI approaches and CTH models
 - Ransomware research directions for detection
 - Available ransomware datasets





Content

I. INTRODUCTION II. BACKGROUND **III. LITERATURE REVIEW** A. CTI STUDIES **B. CTH STUDIES IV. CYBER THREAT INTELLIGENCE TECHNIQUES** V. MALWARE ANALYSIS A. STATIC ANALYSIS **B. DYNAMIC ANALYSIS** C. HYBRID ANALYSIS **VI. CYBER THREAT HUNTING TECHNIQUES** A. TRADITIONAL MACHINE LEARNING APPROACHES **B. DEEP LEARNING APPROACHES** C. OTHER APPROACHES VII. DIRECTION OF FUTURE RESEARCH ON RANSOMWARE VIII. RANSOMWARE DATASETS IX. CONCLUSION





I. INTRODUCTION

- Ransomware **encrypts** a user's files or **locks** the system [3].
- 2021, 66% of surveyed companies were attacked by ransomware [2].
- 1989 (J. Popp), 1st ransomware (AIDS / PCCyborg). Floppy-disks w/ AIDS researches & malicious scripts [4].
- Ransomware attacks have evolved using different tactics and techniques: **Cryptoransomware** (RSA or AES) and **locker**-ransomware [50].
- Attackers primarily design ransomware attacks for **money extortion** from the victims.
- Famous ransomware families include:
 - **CryptoWall**, 2014, spread by phishing email, exploit kits, infected attachments [5].
 - **TeslaCrypt**, 2015, exploit kits [5].
 - Locky, 2016, embedded macros with Microsoft Office documents [5].
 - **Cerber**, 2015, exploit kits, ransomware as a service (RaaS) [5], and
 - **WannaCry**, 2017, SMB vulnerability in MSWindows [6], [7].
- Target platforms: PCs, **mobile devices**, and Internet of Things (**IoT**) **devices** [8].





I. INTRODUCTION

- Most of ransomware **solutions are reactive**: File hashes, IP address, DNS record [9].
- **Proactive approach**: registry path, system calls, DNS queries, run-time activities [11].
- Cyber Threat Hunting (**CTH**) is proactive approach utilized to secure critical assets [12].
- CTH identifies hidden threats, disables them, and sets policies for future-avoidance.
- Cyber Threat Intelligence (CTI) collects information beyond what is available [13]. Evidence-based knowledge outside security logs is necessary to adopt a proactivity.
- The study does a literature review of CTI and CTH for both malware and ransomware.
- Paper organization:
 - a. Section 2 summarizes ransomware studies and existing CTI and CTH techniques.
 - b. Section 3 provides an overview of cyber threat intelligence techniques.
 - c. Section 4 presents a detailed overview of malware analysis approaches.
 - d. Section 5 discusses cyber threat hunting techniques.
 - e. Section 6 discusses the evolution of ransomware attacks and research directions.
 - f. Section 7 discusses datasets of ransomware detection studies.
 - g. Section 8 provides the conclusion of this study.





- Ransomware targets PCs, mobile, cloud-based, IoT, ICS, and other systems [14], [15].
- Researchers have developed **taxonomies** to understand **how ransomware operates**.
- Specific countermeasures to secure different digital assets.
- Ransomware is classified into **two categories** based on confiscated resources:
 - In a Locker-ransomware attack, the victim will not be able to reach system services; however, data will not be compromised. Locker-ransomware is classified by the type of non-data resources it encrypts, such as operating systems, applications, services, user interfaces, and other utilities.
 - Crypto-ransomware encrypts data resources and requests a ransom payment from users. Crypto-ransomware is classified into three types based on the encryption process: symmetric (AES), asymmetric (RSA), and hybrid.
- A deep understanding of the attack steps is required to discover an effective solution.





Most common ransomware attack phases:

- Infection phase: The malicious ransomware code enters the victim's system. Different infection vectors such affiliate programs, exploit kits, and email-based malvertising campaigns [16].
- **Installation phase:** The ransomware installs itself on system and takes control without attracting attention.
- **Communication phase:** Ransomware establishes initial connection with the main adversary in the command and control (C&C) server to carry out following level of actions.
- **Execution phase:** Ransomware begins malicious operations on the victim's resources such encrypting data, deleting files, accessing file systems, locking procedures, modifying master boot records (MBRs) [17].
- **Extortion phase:** Ransomware notifies users and provides instructions.
- **Emancipation phase:** After ransom payment, attackers would send a link to infected victims `that contains a specific decryption tool for some crypto-ransomware attacks.





- Security researchers have investigated two **defense approaches** for ransomware attacks: **signature-based** and **behavior-based** approaches.
 - Signature-based methods, often known as static analysis, refer to the process of examining a malicious file without its execution. Because of the growth of ransomware attacks and anti-forensic tactics such as packing and obfuscation, signature-based approaches have limitations.
 - Behavior-based approaches, often known as dynamic analysis, refer to running a malicious program and observing its activities in the system. Behavior-based approaches can strive for the detailed characteristics of ransomware behavior. Their ability to strive for detailed characteristics makes using a defensive technique based on ransomware behavior much more effective.
- Employing a **behavior-based** approach as a defensive strategy is **more effective** in preventing ransomware attacks.







Figure 1. Steps involved in a typical ransomware attack.





A. CTI STUDIES

- CTI is a proactive method that gathers valuable information from various sources to provide insight into the most recent cyber vulnerabilities and threats.
- Williams at [18] utilized a web crawling technique to find proactive cyber threat intelligence (CTI) approaches in hacker forums. They implemented the Depth-First Search (DFS) technique, an incremental crawling method for collecting attachments while avoiding various popular anti-crawling measures.
- Li at [19] relied on 131 articles focusing on security event-related topics to build an SVM ML-based proactive CTI.

Ref*	Data source	Analytics	Selected threats	
[18]	10 hacker forums	LSTM	mobile, database, Web, sys- tem, network.	
[19]	131 articles	SVM	Event-based security articles.	
[20]	8 forums	SVM and LDA	DDoS, SQLi, Keyloggers, web exploits.	
[21]	9 darknet markets	LSTM	Financial market threats.	
TABLE 1. A summary of CTI techniques				



Seoul Tech II

Ubiquitous Computing & Security Laboratory

A. CTI STUDIES

- **Samtani** at [20] presented a methodology for implementing a proactive CTI by mining hacker communities for source codes, tutorials, and attachments. The framework employs social network analysis methodologies and metrics to identify the key individuals behind discovered hacking assets.
- **Ebrahimi** at [21] focused on cyber threats hosted by the deep net market to avoid significant financial losses. They created a web crawler that used many approaches to combat deep net marketplace anti-crawling mechanisms.

Ref*	Data source	Analytics	Selected threats	
[18]	10 hacker forums	LSTM	mobile, database, Web, sys- tem, network.	
[19]	131 articles	SVM	Event-based security articles.	
[20]	8 forums	SVM and LDA	DDoS, SQLi, Keyloggers, web exploits.	
[21]	9 darknet markets	LSTM	Financial market threats.	
TABLE 1. A summary of CTI techniques				



Seoul Tech

Ubiquitous Computing & Security Laboratory

- Cyber threat hunting (CTH) is an approach that integrates CTI with data analysis methods to detect and respond to threats proactively.
- [23] Homayoun developed sequential pattern mining as a ransomware hunting mechanism.
- [24] Mavroeidis suggested a Sysmon log-based automated threat hunting system.
- [25] Darabian et al. developed an integrated multi-view learning approach that uses multiple features rather than a single feature view. The proposed solution was employed on Windows, Android, and IoT platforms.
- [26] Naik developed triaging methods as hunting techniques to determine the similarity of two ransomware samples. They applied four evaluation methods.





- [27] Jadidi proposed an industrial control system threat hunting framework. The proposed framework focuses on detecting cyber threats.
- [28] Haddad Pajouh developed an IoT malware hunting method using a Long Short-Term Memory (LSTM) structure.
- [29] Jahromi developed an Extreme Learning Machine (ELM) approach that includes two hidden layers.
- [30] Homayoun developed a system for deep ransomware threat hunting in the fog layer. They used LSTM and CNN for classification to discover ransomware attacks within the first 10 seconds of program execution.





- [31] Al-rimy proposed two novel techniques, incremental bagging (iBagging) and enhanced semi-random subspace selection (ESRS) iBagging technique is used to build incremental subsets that show the evolution of crypto-ransomware behavior over various attack phases. ESRS technique is used to construct feature spaces and exclude weak features.
- [32] Al-rimy proposed the Redundancy Coefficient Gradual Upweighting (RCGU) technique which improves redundancy-relevancy tradeoffs during feature selection. RCGU technique increases the redundancy term weight proportional to the number of selected features.
- [33] Kok proposed a Pre-Encryption Detection Algorithm (PEDA) that aims to discover crypto-ransomware attacks at the phase of pre-encryption. The first level uses static analysis to compare the file signature with the known ransomware signature.





- [34] Darem proposed an adaptive behavioral-based incremental batch learning malware variants detection model.
- [35] Roy proposed a deep learning-based ransomware detector (DeepRans) The proposed model monitors the infected host's suspicious activity.
- [36] Pundir proposed a hardware-assisted ransomware detection technique using DL methods. They monitored micro-architectural events using a hardware performance counter to detect abnormal events.
- [37] Ullah proposed a ransomware detection model using online ML classifiers. The model performs detection by tracing the ransomware behavior features during the execution.





- [38] Zhang proposed a deep learning-based model that uses a self-attention mechanism. They used an N-gram of opcodes to identify ransomware fingerprints.
- [39] Khan proposed a DNAact-Ran system that uses digital DNA sequencing along with ML to detect ransomware.
- [40] Poudyal proposed an AI-based ransomware detection framework (AIRaD). The proposed framework combines static and dynamic analysis to detect ransomware attacks.
- [41] Zuhair proposed a machine learning-based multi-layer ransomware detection system. The proposed solution consists of analysis, learning, and detection phases.
- [42] Alrazib proposed DL-driven software-defined networking (SDN) intrusion detection system. They studied the presence of emerging cyber threats in the IoT environment.
- [43] Javeed proposed a novel SDN-enabled hybrid DL-driven cyber threat detection. The proposed solution detects cyber threats on the IoT platform.





B. CTH STUDIES

Limitations and gaps in current CTH solutions are as follows:

- **Disregarding the evolving nature of ransomware or malware attacks:** previous studies focused on **detecting malware** based on **traditional techniques** that compare one or more static features. This is insufficient to the evolving nature of attacks that utilize elusion and offense techniques.
- Relying on using classification methods based on static features: previous works relied on analyzing only static information that is ineffective against sophisticated malware. Static features involve analyzing a malware binary file without executing the code, such as determining the malware's signature and calculating the hash of the malware file.
- Small number of samples in the used dataset: The growth in the complexity of ransomware or malware attacks requires the utilization of a large, diverse, and up-to-date dataset. Some studies used a small dataset that could affect the prediction and lead to overfitting issues.





B. CTH STUDIES

Limitations and gaps in current CTH solutions are as follows (continuation):

- **Imbalanced data:** Studies include classification data with skewed class proportions that make one class a majority class and another a minority class.
- **Hiding performance results:** Performance results, such as model accuracy, f-measure, and other measures, could help other researchers evaluate previous works and solve the current challenges. Some studies do not show performance results and findings that could affect the research field.
- Some studies have not specified the source of the collected data samples: Datasets are a significant part of scientific research. Some previous studies have not described the source of malware or ransomware samples.
- Some studies have focused on finding a set of features that cannot be shown in other versions of ransomware samples.





IV. CYBER THREAT INTELLIGENCE TECHNIQUES

- CTI can provide detailed information related to anticipated cyber attacks. For example, an **email designed for phishing** attacks could include various vital **features** such as the **attack technique** used, **attacker information**, **target information**, **software**, and **tools** used to launch the attack [46].
- The **collection** and **analysis** of massive amounts of online sources of threat data present a new area of challenges that enhance CTI abilities to mitigate or disable rising attacks [20]. To discover online sources, extensive data analysis, awareness of web crawling and anti-crawling mechanisms, understanding of foreign languages, knowledge of cyber world terms, and understanding of the complex structures of malicious assets are needed.
- The **web crawling mechanism** is applied to search for web content. A web crawler is used for different purposes, such as **searching for and extracting information** or classifying web content. A crawler **parses HTML tags** and **retrieves pages**, **extracts new hyperlinks** from these tags, and stores HTML content. After collecting the data, the analysis technique is utilized to leverage the discovered information to understand the critical trends of malicious cyber assets.





V. MALWARE ANALYSIS

- Malware analysis could be classified into static, dynamic, and hybrid.
- Malware samples can be analyzed manually or automatically.
- STATIC ANALYSIS:
 - Static analysis: reverse engineering, disassembling, dissecting binary file.
 - The malware structure is identified by static analysis with no code execution.

• DYNAMIC ANALYSIS:

- Dynamic analysis by observing or debugging a malware's program instructions.
- Isolated environments such as virtual machines or sandboxes are used.

• HYBRID ANALYSIS:

- Hybrid analysis is a file analysis with both static and dynamic analysis aspects.
- Static analysis is easier and faster than dynamic analysis.
- Impossible to analyze malicious software that uses obfuscation, packed, or polymorphic techniques using static analysis.
- Dynamic analysis shows malware's actual functionality.





A. TRADITIONAL MACHINE LEARNING APPROACHES

- ML, machines learn from data or experience to automate the building of analytical models.
- Supervised learning, unsupervised, semi-supervised, and reinforcement learning are the four major ML approaches.
- ML involves data collection, cleaning and preparation, model building, model evaluation, and model deployment.





Ref*	Attack	Features	Method	Results
[23]	Ransomware	13 selected features	J48, Random Forest (RF), Bag- ging and MLP	Achieved F-Measure of more than 0.98 with FPR of less than 0.007.
[25]	Malware	OpCodes, ByteCodes, header information, permission, attacker's intent, and API calls	SVM	The accuracy of the proposed model is 99.6% on IoT dataset, 99.6% on Android dataset, and 98.01% on Windows dataset.
[31]	Ransomware	API calls	Ensemble-based learning	The detection accuracy ranges between 0.957838 to 0.97885 based on the number of the selected features.
[32]	Ransomware	API calls	SVM, Logistic Regression, De- cision Tree, KNN, RF, Ad- aBoost, and MLP	The detection accuracy, detection rate and false positive of all feature sets per classifier ranges between 0.9573-0.9909, 0.9621-0.9940, and 0.0384-0.0068.
[33]	Ransomware	API calls	RF	The achieved recall rate was 100% based on 80:20 ratios of training and testing, and 99.9% recall rate with a 10-fold cross-verification test.
[37]	Ransomware	API calls of registry, network, and file system activities.	Modified DT, RF, and Ad- aBoost.	The achieved detection accuracy of Modified DT, RF, and AdaBoost are 99.56%, 99.24%, and 98.37%, respectively.
[39]	Ransomware	Generated DNA sequence of selected features.	Naïve Bayes (NB), RF, and se- quential minimal optimization (SMO).	The achieved detection accuracy is 87.91%.
[40]	Ransomware	Multi-level of static and dy- namic features.	SVM, LR, RF, AdaBoost with J48, and J48.	SVM and AdaBoost with J48 achieved the highest accuracy of 99.54%. J48 classifier achieved the second highest accuracy of 99.26%.
[41]	Ransomware	9 selected features.	DT and NB.	The model achieved an average accuracy of 96.27%.



Tabla 2. A summary of ML methods used for CTH techniques.







B. DEEP LEARNING APPROACHES

- Deep learning (DL) is a subset of ML that learns from data.
- DL models require a large amount of data for each problem domain.
- DL algorithms require high computational capabilities to train models with a large amount of data.





Ref*	Attack	Features	Method	Results
[28]	Malware	OpCodes	LSTM	The proposed model achieved 98% detection accuracy against IoT malware samples not used in model training.
[29]	Malware	OpCode and system calls	TELM	TELM outperformed the original ELM models. Achieved ac- curacy rates are from 95.80% to 99.03%.
[30]	Ransomware	Sequence of actions	LSTM, CNN, and MLP	LSTM outperformed the other methods and achieved 0.996 F- measure of detecting ransomware in binary classification, and 0.972 TPR of identifying ransomware family in multi-class classification.
[34]	Malware	API calls	Sequential deep learning	The achieved F-measure results are higher than 99% and the average accuracy of detecting new malware is 99.41%.
[35]	Ransomware	Event ID of real-time bare metal logs	Attention-based bi-LSTM	Achieved 99.87% detection accuracy and 99.02% F-measure for early detection. Also, the model achieved 96.5% detection accuracy for classifying abnormal events.
[36]	Ransomware	Micro-architectural event traces	RNN and LSTM	Achieved an average of 97% detection accuracy.
[38]	Ransomware	N-gram of OpCode	Self attention-based CNN.	The achieved detection accuracy is 0.895 and the average F-measure is 0.873.
[42]	Cyber	Multiple features.	Hybrid DNN-LSTM	The achieved accuracy is 99.55% and the average F-measure is 99.42%.
[43]	Cyber	Network flow features	Hybrid LSTM-GRU	The achieved detection accuracy is 99.74% and the average F-measure is 99.79%.



Tabla 3. A summary of DL methods used for CTH techniques.



C. OTHER DATA ANALYSIS APPROACHES

• Other data analysis approaches that do not include artificial intelligence methods have been utilized in CTH studies such as, [24], [26], and [27].

Ref*	Attack	Features	Method	Results
[24]	Malware	Features extracted from Sys- mon logs	Dynamic analysis	The performance of the model has not been evaluated.
[26]	Ransomware	API names and order, Byte structure, textual, and binary patterns	IMPHASH, SSDEEP, SDHASH and YARA rules	SDHASH outperformed the other three methods based on the total number of samples matched.
[27]	Malware	events	MITRE ATT&CK matrix and a diamond model of intrusion analysis	The proposed framework was evaluated using three scenarios.

Tabla 4. summary of other data analysis techniques used for CTH.





- Most ransomware-related research works focus on different characteristics such as threat delivery, encryption algorithm and communication, associated IoCs, and behavior analysis [53].
- New ransomware variants are constantly being developed and current solutions that rely on static analysis only detects only earlier forms of ransomware samples. Therefore, cybercriminals just apply advanced techniques to conceal the ransomware execution and avoid detection.
- Ransomware can appear as a standalone crypto-worm that replicates itself to other computers to maximize the impact on the network. In addition, ransomware can appear as a Ransomware-as-a-Service (RaaS), which is a distribution kit sold on the dark web. RaaS permits novel attackers with limited technical skills to launch ransomware attacks [55].
- Ransomware can infect files on locally fixed, removable, or remotely shared drives. To minimize detection, attackers may attempt to sign their ransomware using code-signing technology by buying or stealing it.
- Current ransomware utilizes exploits to elevate their privileges. After that, the ransomware will start encrypting as many files as possible to ensure receiving ransom money from the victim.





- A ransomware attack has a financial impact on an organization when it encrypts its mapped network drives. Restoring multiple servers from backup data takes a long time, and data could not be up to date. Many organizations use only backup solutions as a critical defense against ransomware, which makes it a recovery solution rather than a detection solution.
- The ransomware performs the file encryption process using two methods: overwrite (in-place) and copy. The overwrite method encrypts files by reading the original file, writing an encrypted version over the original file, and renaming the file. On the other hand, the copy method encrypts files by reading the original file, creating an encrypted copy, and deleting the original file. It is impossible to recover the original files using the overwrite method. However, ransomware that uses the copy method will use an additional wiping action to ensure that data files are not recoverable.
- Ransomware behavior follows specific patterns that include the file identification process, file encryption, network command, and control communications [56]. In most ways, ransomware uses a Windows application programming interface (API) to make function calls. Windows API offers a collection of programming interfaces that simplify the software development process. Windows API calls can be used as behavioral features to identify abnormal patterns.





28

#	Category	Windows API call	Purpose
1	1 registry RegOpenKey		Open a registry key for editing and querying.
		RegCreateKey	Create the specified registry key.
		RegQueryValue	Retrieve the type and data for the specified value name.
		RegSetValue	Add a new value to the registry & sets its data.
		RegEnumValue	Enumerate the values for the specified open registry key.
2	fileSystem	CreateFile	Open document for reading and writing.
		ReadFile	Read data from document.
		WriteFile	Write data into document.
3	system	LoadLibrary	Load the specified module into the address space.
4	process	CreateRemoteThread	Create a thread of another process.
		CreateProcessInternal	Create a new process and its primary thread.
		ShellExecute	Perform an operation on a specified file.
		ExitProcess	End the calling process and all its threads.
5	memory	VirtualAlloc	Reserve, commit, or change the state of a region of memory within the virtual
			address space of a specified process.
6 synchronization CreateMutex Create or opens		CreateMutex	Create or opens a named or unnamed mutex object.
		OpenMutex	Get a handle to another process's mutex.
7	services	OpenSCManager	Return a handle to the Service Control Manager.
OpenService Open an existing service.		Open an existing service.	

Table 5. Windows API calls categories and examples extracted from different ransomware samples.





29



Figure 3. Extraction of API calls sequence from ransomware sample.





VIII. RANSOMWARE DATASETS

- Datasets are essential to foster the development of an effective ransomware detection solution.
- Therefore, the accuracy of the solution is directly related to and dependent on the input dataset.
- Datasets contain several samples for benign and ransomware.
- Datasets for ML could either be privately collected or publicly available to anyone.
- Different ransomware studies used datasets from different repositories.
- Popular repositories that offer malware data include the following sources: VirusTotal, VirusShare, and the Zoo.





VIII. RANSOMWARE DATASETS

Ref*	Attack	Platform	Dataset classification	Dataset source	Number of samples
[23]	Ransomware	Windows	Ransomware	virustotal.com	1624
			Benign	portableapps.com	220
[26]	Ransomware	-	ransomware	hybrid-analysis.com malshare.com	200
[30]	Ransomware	Windows	Ransomware	virustotal.com	660
			Benign	portableapps.com	219
[31]	Ransomware	Windows	Ransomware	virustotal.com	8152
			Benign	informer.com	1000
[32]	Ransomware	Windows	Ransomware	virusshare.com	39378
			Benign	informer.com	16057
[33]	Ransomware	Windows	Ransomware	Sgandurra [58]	582
			Benign		942
[35]	Ransomware	Windows	Ransomware	PC host logs	4,820
			Benign		1,033,297
[36]	Ransomware	Windows	Ransomware	virusshare.com	80
			Benign	OpenSSC C programs	76
[37]	Ransomware	Windows	Ransomware	virustotal.com	35369
			Benign	_	43191
[38]	Ransomware	Windows	Ransomware	virustotal.com	1787
1000			Benign	System executable files	100
[39]	Ransomware	Windows	Ransomware	github	582
			Benign	System executable files	942
[40]	Ransomware	Windows	Ransomware	virustotal.com	550
			Benign	Windows 10 open-source software	540
[41]	Ransomware	Windows	Ransomware	virustotal.com malware blacklist	10000
- 1939 - Di			Benign	collected manually from websites	500

Tabla 6. Datasets used by state-of-the-art ransomware detection works.





IX. CONCLUSION

- Ransomware is an evolving form of malware designed to block access to the system or encrypt its data.
- Various static and dynamic features of ransomware can be extracted and used to reveal its activities.
- This paper presents a systematic review of Cyber Threat hunting techniques for detecting ransomware attacks.
- The previous works of CTI and CTH have been investigated, and the limitations and gaps have been mentioned.
- Then, we explained the CTI technique. We provided an extensive overview of the malware analysis.
- CTH techniques are discussed based on the used data analysis method.
- Ransomware evolution and research directions are highlighted.
- The available ransomware datasets used in the previous works are mentioned with their data sources.





IX. CONCLUSION

- In summary, ransomware attacks must be detected proactively, as shown in this study.
- Developing an effective ransomware CTH technique that can detect known and unknown ransomware is a concern.
- We provided a detailed review of ransomware research directions and the available ransomware datasets utilized with different data analysis methods.
- In our future work, we will adopt a CTI method to enhance the development of a CTH technique by collecting the latest shared information about ransomware attacks.
- Subsequently, the collected information will be incorporated into an effective new learning strategy model to enhance detection accuracy.
- A deep focus on dynamic features will be performed to hunt ransomware attacks based on behavior classification.



