# Holistic Cyber Threat Hunting Using Network Traffic Intrusion Detection Analysis for Ransomware Attacks

**Kanti Singh Sangher, Arti Noor, and V. K. Sharma**

**Abstract** In recent times, cybercriminals have penetrated diverse areas or sectors of the human business enterprise to initiate ransomware attacks against information technology infrastructure. They demand for money called ransom from organizations and individuals to save valuable data. There are varieties of ransomware attacks floating worldwide using intelligent algorithms and with the usage of different setup vulnerabilities. In our research work, we are exploring the latest trends in terms of sector-wise infiltration, captured the most popular among available and also the distribution of the number of attacks using the location information available at the country level. To achieve the correlation between the sectors and locations along with the parametric analysis, we have utilized artificial intelligence techniques. Accuracy of the prediction of attack based on the sector level analysis we have implemented Random Forest and XGBoost algorithm. This research work focuses primarily on two aspects, first is to explore the different aspects of ransomware attacks using intelligent machine learning algorithms. The method used insights to severity of spread of ransomware attacks, second research outcome is to forensically evidence finding of the attack traces using traffic analysis. The challenge is to learn from the previous weaknesses available in the infrastructure and at the same time to prepare the organization and countries' own prevention methods based on the lessons learnt, our exploratory analysis using the latest set of data implementing with AI will give a positive dimension in this area. Also, the proactive approach for managing the data safely is based on the finding of digital forensic analysis of infected ransomware traffic.

K. S. Sangher (✉) · A. Noor · V. K. Sharma
School of IT, Centre for Development of Advanced Computing, Noida 201307, India
e-mail: kantisingh@cdac.in

A. Noor
e-mail: artinoor@cdac.in

V. K. Sharma
e-mail: vksharma@cdac.in

## 1 Introduction

Ransomware is one of the most sophisticated online threats, dominating security exploitation at both the individual and organizational levels. Data loss is unaffordable because of the sensitivity tied to it and variants are becoming more harmful as time goes on. Therefore, it is imperative to conduct intelligence analysis in order to filter the most recent attacks [1]. If we can manage that, it will enable us to improve the infrastructure, fortify the environment, and, most critically, be well prepared to thwart future attacks.

The security dangers are changing as well as our dependence on digital technology grows in both our personal and professional lives. One of the prominent threats is malware heavily damaging the cyberspace. Ransomware and its impact have changed the reach of malware in worldwide, once infected the resource device restricts the access of files and folders till the ransom raised by the cybercriminal paid, mostly nowadays in digital currency such as Bitcoin to get the data back. Recent trends show that ransomware is not limited to a particular domain but penetrates different sectors such as education, health, information technology, business, and research. Nature of the attack and its consequences are tricky in the case of ransomware as the damage is mostly irreversible even after the removal of the malware that caused the attack. Hence, cyber security becomes a critical concern for researchers and organizations to find the solution to overcome ransomware attacks or to be prepared with preventive solutions. Recently, ransomware has matured in intricacy, difficulty, and diversity to turn into the most vicious among the existing malware trends. In addition to this, Cisco's annual security reported that ransomware is rising at a yearly rate of over 300%. The method used gives insights to severity of spread of ransomware attacks, and the second research outcome is to forensically finding the evidence of the attack traces using traffic analysis [2].

## 2 Literature Survey

Critical infrastructure is severely impacted by malware, sometimes known as malicious software. The goal of these is to harm the victim's computer or service networks. Malware comes in many different forms, including viruses, ransomware, and spyware. Malware known as ransomware has been shown to use complex attack methods that have undergone numerous mutations. Many different sectors of the economy, including transportation, telecommunications, finance, public safety, and

health services, have been impacted by ransomware. User data is made unavailable or unreachable using crypto modules incorporated in malware. The organization is required to pay a ransom to regain access once ransomware either locks the equipment or encrypts the files. With code obfuscation, various ransomware display polymorphic and metamorphic tendencies, making it difficult to scan for and identify them with current methods. The market is classified by deployment, application, location, and other factors for the many open-source and commercial anti-ransomware programs. AOKasperskyLab, FireEyeInc, MalwarebytesInc, SophosLtd, SymantecCorporation, SentinelOneInc, ZscalerInc, TrendMicro Incorporated, and more top companies offer ransomware prevention solutions. Using software libraries and sandboxes like CryptoHunt, Cryptosearcher, etc., ransomware is discovered by conducting a crypto module search and analysis.

The new ransomware variants are not recognized by the existing solution, and the effects are only discovered after the attack, despite the fact that there are ongoing upgrades or improvements to the existing anti-ransomware systems. Analysis also demonstrates that no ransomware variant changes after each infection, allowing ransomware authors to stay one step ahead of their victims because the ransomware's associated signatures, domains, and IP addresses become dated and are no longer recognizable by threat intelligence and signature-based security tools. In the critical infrastructure sectors, the requirement for a new paradigm to recognize the new and evolved ransomware is crucial. These ransomware attacks have the potential to cause significant financial harm and substantial losses, making cyber-insurances for all enterprises necessary.

The majority of ransomware detection programs use behavioral detection, often known as "dynamic analysis" [3–6] based on dynamic analysis using an SVM classifier. They initially retrieved a specific ransomware component known as the Application Programming Interface (API) call, after which they used Cuckoo Sandbox to analyze the API call history and its behaviur. Q-gram vectors serve as a representation for the API calls. They employed 312 goodware files and 276 ransomware files. The findings showed that using SVM, ransomware could be detected with an accuracy of 97.48%. Vinayakumar et al.'s [5] new approach suggested collecting the API sequences from a sandbox utilizing dynamic analysis.

## 3 Present Situation of Ransomware Worldwide and India's Stack

India is emerging as a cyber power in the international community and at the same time by the end of this year, about 60% of the Indian population (840 million), will have access to the internet, claims a report. The flip side of this is, increasing cybercrimes. Over 18 million cases of cyber attacks and threats were recorded within the first three months of 2022 in India, with an average of nearly 200,000 threats

every day, according to the cyber security firm Norton. So, there is a strong need to be more vigilant, proactive, and smart while handling the cybercrimes.

## 3.1  Finding Recent Heavily Used Set of Ransomware Attacks and Their Behavior

The attackers are using the more complex, destructive, and easy-to-execute ransomware variants [7]. The dataset used in our research work provided by the DSCI sources consists of the worldwide penetration along with the per sector as target. The dataset "Ransomware Attacks" has the following attributes: description, sector, organization size, revenue, cost, ransom cost, data note, ransom paid, YEAR, YEAR code, month, location, Ransomware, no of employees Source Name, URL details, etc. The training environment for the ransomware attack dataset is set up using Anaconda v3. Environment implemented using IPython utilizing Jupyter Notebook from Google Collab. The word cloud created from the dataset indicates that REvil, Cryptowall, and Ryuk are some of the most recent trends. Even a small number of the recent ransomware attacks have unclear sources (Fig. 1).

For different operating systems and gadgets, there are many ways that ransomware actually gets onto the computer. One of the most notable assaults in 2021 was REvil, which successfully hacked more than 1500 supply chain management-based firms and spread unintentionally to all associated clients. A remote IT management service provider, Kaseya MSP, was one of those compromised and used to distribute the ransomware REvil/Sodinokibi. The ransom demand was made by evil threat actors, who demanded between $50 million and $70 million in Bitcoin to unlock all the encrypted information. Virtual Systems Administrator (Kasseya VSA) is a remote management and monitoring application for networks and endpoints used by businesses and MSPs. The Kaseya network was breached by the REvil attackers, who then took control of the VSA software update system to deploy the initial payload via the Kaseya agent.

*C:\Program Files (x86)/Kaseya/{ID}/AgentMon.exe*

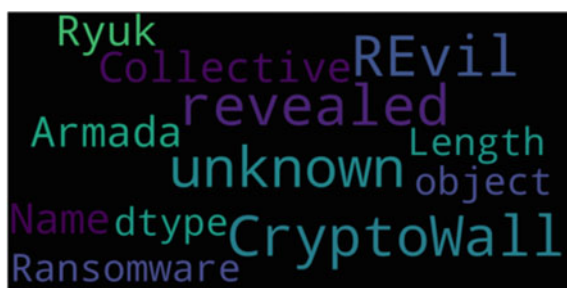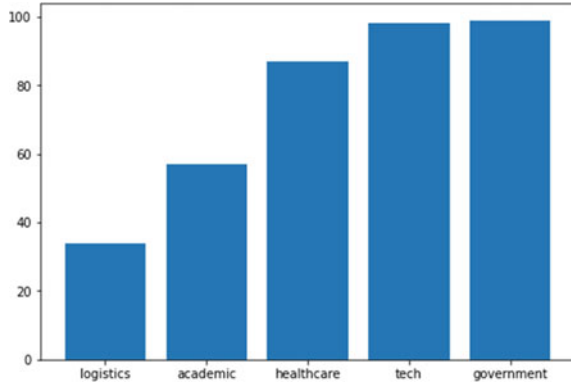**Fig. 1** Word cloud created using the ransomware attribute from the dataset

**Fig. 2** Existing and emerging targets for ransomware attacks till 2021



The launch method used to intrude initiates payload by following the precaution to get detected. Again to make things perfect it removes the certutil.exe (*cert.exe*) and *agent.crt* transitional binaries too. The first threat intrusion seems genuine, being just an older version of Windows Defender binary (*msmpeng.exe*). The next one is a binary called REvil encryptor dll (*mpsvc.dll*). This facilitates side-loading of *mpsvc.dll* into *msmpeng.exe* to access the data.

Cybercriminals have understood that supply chain attacks are the new area to explore and make profit. Various MSPs or IT infrastructure management platforms are using agent software to provide organizations different services. If the security is breached of these deployments then they will be used to distribute harmful malware. It is a high time to protect the different domains which utilize the network and remote machines to share services the countermeasures and spread within the IT community should be availed [8]. So, based on the recent attacks dataset shared from the Data Security Council of India (DSCI) a premier industry body on data protection in India, setup by NASSCOM. We have used the AI techniques and performed analysis to gather target organizations details for ransomware attacks in last 5 years. Figure 2 shows the outcome from the analysis of the dataset in the form of a graph, which precisely depicts that, in recent time, government organizations are the prime target along with the technology-based industries, and then health care is identified as an attractive spot for cybercriminals where a lot of personal health records and medicinal data are floating further logistics is the upcoming area which is facing lot of variants of ransomware attacks [9].

## 3.2 Intelligent Discovery of the Attack Trends

Using the experiments on the dataset we tried to visualize the attack reach at the national level with a comparison to worldwide data. It clearly shows that India is one of the major targets and govt. organizations with technology are sharing the target. So, there is an immediate need to strengthen our govt. organizations to protect the
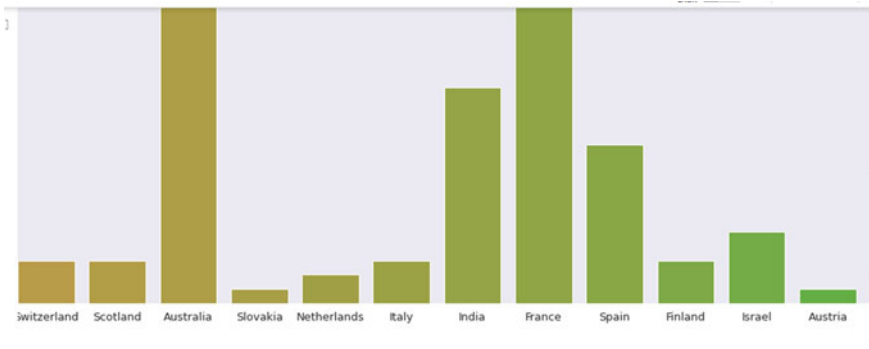
**Fig. 3** India's stack in worldwide statistics in ransomware attacks till 2021

environment to safeguard the data due to its sensitivity and the damage it can cause for our nation (Fig. 3).

One of the significant analyses shows the correlation between the different parameters of the dataset, which gives insights to understand the behavioral pattern after the incident happens as we observed from our intelligent analysis using the AI environment that per year cost of ransom is increasing [10] not only in terms of digital transaction but also due to the inclusion of cryptocurrencies, even is some incident or stories which are considered as case studies to prepare the best case practices also experience that after the ransom amount paid by the victim organization cybercriminals not sharing the decryption key on top of that shared decryption key not able to recover the complete data. Figure 4 depicts the correlation in years with respect to ransom cost.

After finding the latest trends of the attack in terms of most common ransomware attacks and the organization level filtration visualization. The next research work
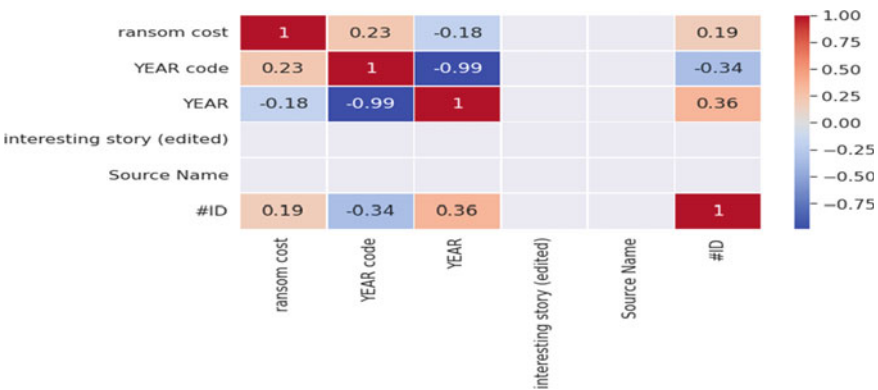


**Fig. 4** Correlation which depicts the YEARS and cost of ransom rose

done was to perform the machine learning algorithm implementation for prediction of precise output (in terms of sector wise/organization calibration).

The correlation matrix's cells represent the "correlation coefficient" between the two variables corresponding to the cell's row and column. A strong positive relationship between X and Y is indicated by values that are close to $+1$, while a strong negative relationship between X and Y is shown by values that are close to $-1$. Values close to 0 indicate that there is no link at all between X and Y. The ransom price in this case varies without regularly growing or decreasing and is unpredictable. Numerous assaults have occurred recently, and the ransom price is high yet inconsistent. Additionally, the type of attack vector and the victims' response affect the total ransom price for ransomware attacks that is paid on an annual basis. This research work will provide very important inputs to the national-level security application team to check and prepare their environment [11], if the victim is aware, then the prevention will be better. After the training of the dataset, it splits the dataset and categorizes it into training and testing sets. The training set includes the classification information and based on the test data we implemented three algorithms, namely, Random Forest, XGBoost, and KNN. A composite data repository is produced by extracting all the features and combining them with the dataset of ransomware attacks. When the composite dataset was prepared, we ran two phases of tests on the same data repository. To determine the accuracy level in the first stage, we used the Random Forest and XGBoost algorithms, two machine learning methodologies. We also noted how long both algorithms took to process.

XGBoost outperforms the competition because it is forced to learn from its mistakes made in the prior stage (so does every boosting algorithm). The main drawback to using XGBoost is if you have a lot of categorical data. Then, you must perform One Hot Encoding, which results in the creation of more features. The overall line is that you can utilize XGBoost for better outcomes, but be sure to perform the proper preprocessing (Keep an eye on datatypes of features). To avoid overfitting, they employ the random subspace approach and bagging. A random forest can readily handle missing data if they are implemented properly. In the second phase, we used K-Nearest Neighbor (KNN) and assessed its accuracy and processing speed using the same dataset. By making calls to the local repository of Jupyter Notebook, we imported the required packages and libraries. The outcomes of applying various machine learning techniques to the ransomware dataset are displayed in the following screenshots. The results found shown in Fig. 5, indicates that an accuracy of implied machine learning algorithm, i.e., 93%, 92%, and 94%, respectively, for Random Forest, KGBoost, and KNN model.

For machine learning when a model fits the training set of data too closely, it is said to be overfit, and as a result, it cannot make reliable predictions on the test set of data. This means that the model has only memorized certain patterns and noise in the training data and is not adaptable enough to make predictions on actual data. However, recall was set to 1 for the purposes of our research, and the data reports we received had varying degrees of accuracy. In order to find ransomware, Kharraz et al. [3] used a dynamic analytic method named UNVEIL. In order to find ransomware, the system builds a fake yet realistic execution environment. About

```
pieus=Classifier.predict(x_test)
from sklearn.metrics import accuracy_score
import xgboost as xgb

xgb=xgb.XGBClassifier()
xgb.fit(X_train,y_train)
preds2=xgb.predict(X_test)
xgb_accuracy=accuracy_score(preds2,y_test)
rf_accuracy=accuracy_score(preds,y_test)
print("Random Forest Model accuracy",rf_accuracy)
print("XGBoost Model accuracy",xgb_accuracy)
knn = KNeighborsClassifier(n_neighbors=1, metric='euclidean')
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
knn_accuracy=accuracy_score(y_pred,y_test)
print("KNN Model accuracy",knn_accuracy)
```

**Fig. 5** Machine learning implementation

96.3% of the time, this system was accurate. Sequential Pattern Mining was employed as a candidate feature to be used as input to the machine learning algorithms (MLP, Bagging, Random Forest, and J48) for classification purposes in the framework ransomware detection system proposed by Homayoun et al. [6].

As a result, the KNN model produced the greatest results in this investigation. Simply said, KNN only stores a footprint of the training data within the model and doesn't really perform any training. KNN's logic is found in the predict() call of its inference step, which is where it uses previously provided training data to identify the k nearest neighbors for the newly supplied instance and predicts the label. For small- to medium-sized datasets, KNN is probably faster than the majority of other models.

## 4 Experimental Analysis of Incident

In our research work, the experimental analysis of ransomware packet capture was done using the Wireshark tool. It shows the step-by-step analysis to find out the traces of intrusion injected to perform the attack. A .tar file of WannaCry has been used as a source for analysis, as it depicts the causes or vulnerabilities within the system that allowed the successful execution of the attack; we propose the defense mechanism to protect the asset within the organization's/individual's infrastructure.

### 4.1 Infected Network Traffic Analysis

The malware's ability to remotely access files is what started the attack, and the attackers used a variety of mechanisms to carry it out. The Server Message Block (SMB) protocol, which is mainly used for file sharing, printing services, and communication between computers on a network, was abused in this case. Ransomware like

WannaCry took use of SMBv1 vulnerabilities by using them. All of these exploits go by the term "Eternal" X. EternalBlue is the one with which most people are familiar. EternalBlue was developed because SMBv1 cannot handle specific packets produced by a remote attacker, which can lead to remote code execution. The WannaCry virus reportedly infected over 230k computers in 2017 and other malware caused over $1 billion in losses globally.

When WannaCry is installed on the network, the following significant things take place:

- The development of files with the WannaCry document extension specifically for encrypting files.
- Ports TCP 445 and 139 of SMBv1's outgoing communication.
- The domain's DNS requests for iuqerfsodp9ifjaposdfjhgosurijfaewrwergwea.com.
- New entries into Windows registry.

All four of these events are trackable and observable. Finding malware that spreads across networks, like WannaCry, requires monitoring. Before the ".WNCRY" extension is introduced, WannaCry encrypts several distinct types of documents. The ransomware connects to the IPC$ share on the remote system after the initial SMB handshake, which includes a protocol negotiation request/response and a session setup request/response, is shown in this study paper through analysis of the network data (Fig. 6).

The malware's ability to connect to a hardcoded local IP is another related component of this attack. The majority of ransomware attacks, as was already noted, use DNS tunneling to create both bidirectional and unidirectional communication between an attacker and the systems on your network. The threat actor can hide until their attack is almost complete if the DNS action is not safe. The system was
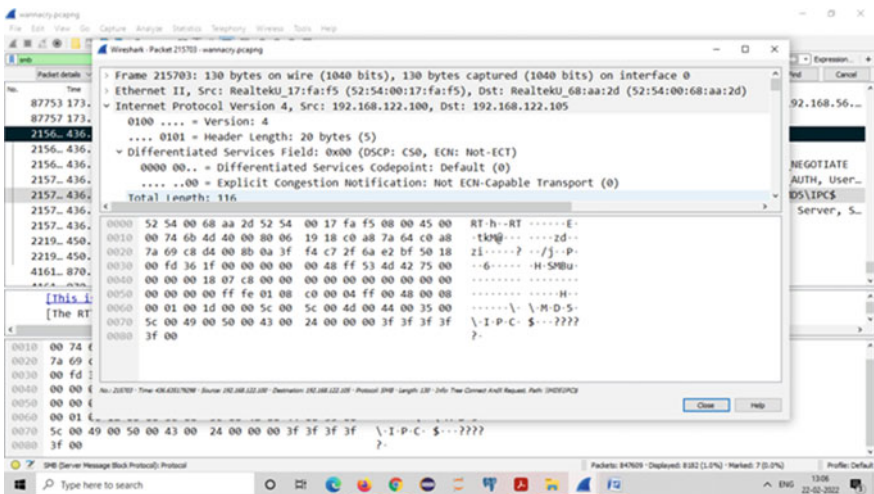


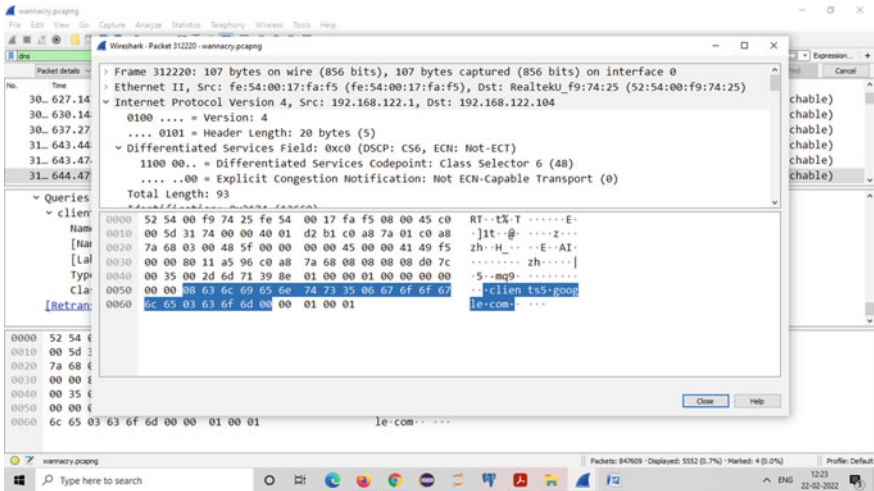**Fig. 6** IPC$ share on the remote machine

**Fig. 7** Malware in traffic

vulnerable to ransomware execution since malware was present, which was discovered during the DNS traffic filtering process. Fig. 7 displays traffic-related malware. A malicious website is clientts5.google.com.

It then sends a first NT Trans request with a huge payload size and a string of Nops, as seen in Fig. 8. In essence, it relocates the SMB server state machine to the vulnerability's location so the attacker can make use of it. It then transmits a first NT Trans request with a huge payload size that is made up of a string of NOPs, as seen in Fig. 8.

To enable the attacker to use a specially constructed packet to exploit the vulnerability, it essentially moves the SMB server state machine to the place where it is present. The next step is to check whether the payload has been successfully installed. If it has, then the SMB MULTIPLEX ID = 82 will be found in one of the packets. The same has been done in this experimental analysis using the filter in Wireshark for a stream of packets and shown in SMB MULTIPLEX ID = 82.

The attack was launched utilizing the SRV Driver exploit MS17-010:Buffer EternalBlue's Overflow. In the worst-case scenario, if an attacker sends specially designed messages to a Microsoft Server Message Block 1.0 (SMBv1) server, they could execute remote code (Fig. 9).

One packet, as seen in the screenshot, signifies that the payload has been installed successfully and that the attacker has run remote code on the victim network. The SMB MultiplexID = 82 field is one of the crucial fingerprints for this attack's success. The contents of the packets can be seen by right-clicking on the Trans2 packet and choosing to Follow -> TCP Stream. The contents of the payloads that caused the buffer overflow and sent the payload necessary for this exploit are shown here (Fig. 10).
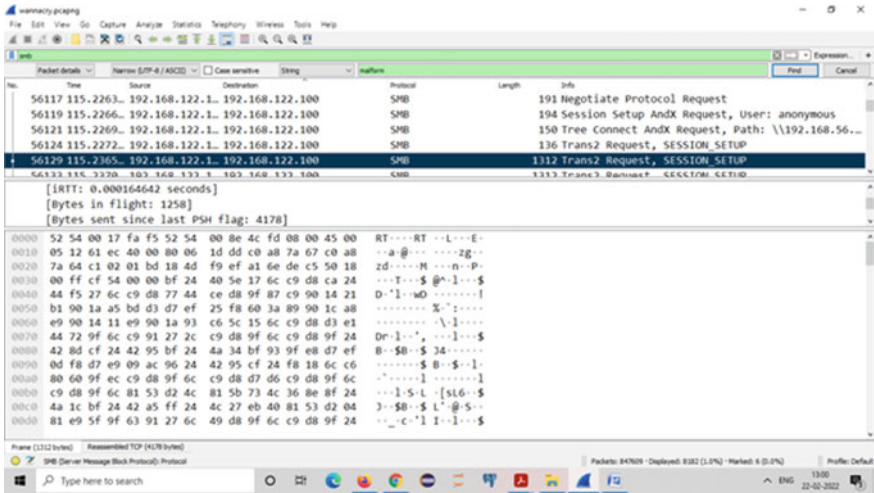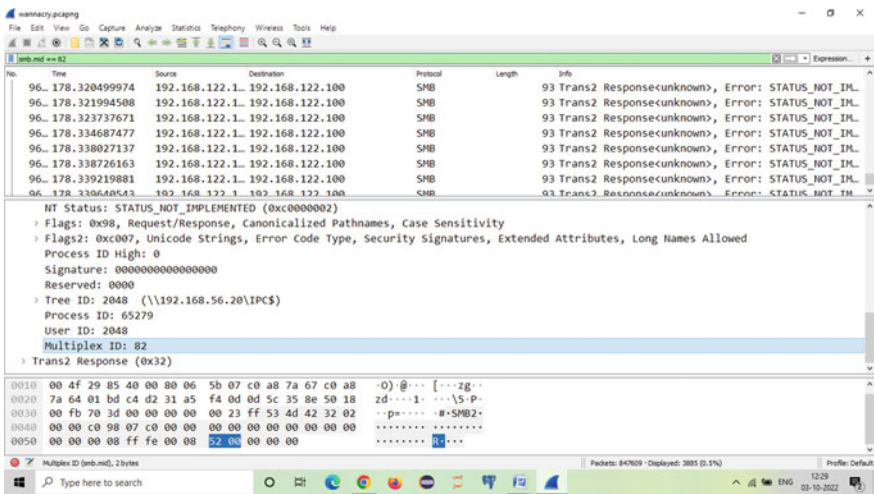
**Fig. 8** NT Trans request with a sizable payload



**Fig. 9** SMB MULTIPLEX ID = 82 within the selected packet

## 4.2 Preventive Measures

The user can recognize this occurrence if there is any folder auditing on Windows folders. The following two entries can also be found in the Windows registry:

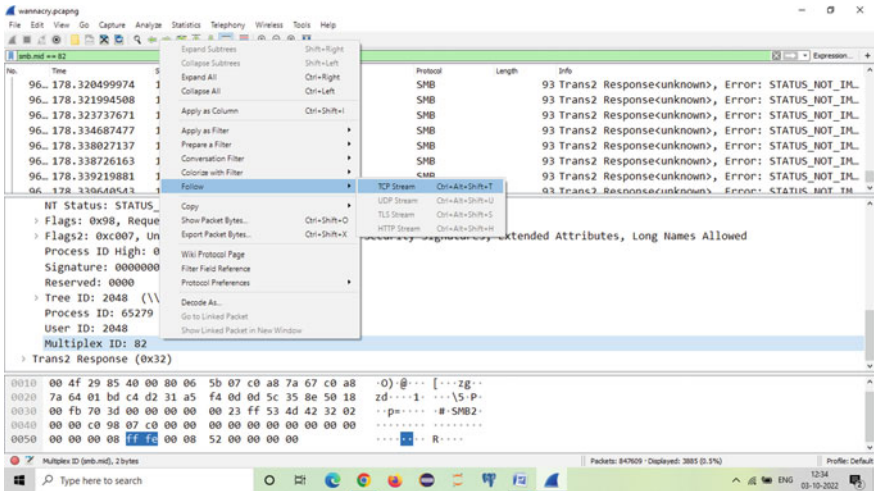- HKLM\SOFTWARE\Microsoft\Windows\CurrentVersion\Run\\<random string> = "<malware working directory>\tasksche.exe"

**Fig. 10** To view the contents of the packets

- HKLM\SOFTWARE\WanaCrypt0r\\wd = "<malware working directory>".

On the local server, users can manually check for file creation and registry key entries or execute a PowerShell script to check for these occurrences. The other two incidents can be investigated using monitoring tools. Additionally, system files must be examined for incoming TCP requests on ports 139 and 445. A rise in requests on these two SMB-specific ports should raise red flags. DNS monitoring and analysis is the second method for finding WannaCry. To protect the system one more check will be helpful that is instigating the browser's unique User_agent value. So, the first thing is to secure the system to prevent ransomware attacks, but if the system security is not strong enough, then learning from the entry points of the incident can help to make the secure environment.

## 5 Conclusion

The paper proposes to intelligently analysis of the ransomware attack data in recent years to visualize the pattern and target of the popular attack, and also how to harden the security to prevent the system from such attacks. This will help tremendously to prepare and harden the organizations' security infrastructure to protect as well as detect the intrusion. The results present intelligent algorithm solutions to ransomware attack penetration at organization level along with the latest set of attacks floating at the worldwide level [12] and also analyze the infected network traffic to find the traces of the ransomware execution using the tool. The research finding gives insightful directions to be aware of the existing threats and prepare the cyber resilience

environment and platforms where targets are identified and well advanced to enable the systems to fight and protect it from the threat vectors.

## 6 Future Scope

The future scope of the present work can be the initial root cause analysis or, in other words, finding the vulnerabilities that were the reason behind the success of cybercriminals. Root cause analysis of the attack will definitely help the organizations' security infrastructure handling team to understand vulnerabilities that helps ransomware to enter in their deployments [13]. Post-incident analysis always gives feedback to be prepared for prevention measures as it helps to serve to plan to handle the situation if the incident happens. Recovery phase also needs to be analyzed due to the uncertainty of the data recovery from the ransomware attacks. A few very basic common points in ransomware attacks are browser exploitation, email, etc. Other vulnerability exploration can be improved in future and applying deep learning for intelligent analysis will be a great area to work.

## References

1. Connolly LY, Wall DS, Lang M, Oddson B (2020) An empirical study of ransomware attacks on organizations: an assessment of severity and salient factors affecting vulnerability. J Cybersecur 6(1):tyaa023. https://doi.org/10.1093/cybsec/tyaa023
2. Dasgupta D, Akhtar Z, Sen S (2022) Machine learning in cybersecurity: a comprehensive survey. J Def Model Simul 19(1):57–106. https://doi.org/10.1177/1548512920951275
3. Lee JK, Chang Y, Kwon HY et al (2020) Reconciliation of privacy with preventive cybersecurity: the bright internet approach. Inf Syst Front 22:45–57. https://doi.org/10.1007/s10796-020-09984-5
4. Kirda E (2017) Unveil: a large-scale, automated approach to detecting ransomware (keynote). In: 2017 IEEE 24th international conference on software analysis, evolution and reengineering (SANER). IEEE Computer Society, pp 1–1
5. Takeuchi Y, Sakai K, Fukumoto S Detecting ransomware using support vector machines. In: Proceedings of the 47th international conference on parallel processing companion (ICPP '18). Association for Computing Machinery, pp 1–6. https://doi.org/10.1145/3229710.3229726
6. Vinayakumar R et al (2017) Evaluating shallow and deep networks for ransomware detection and classification. In: International conference on advances in computing, communications and informatics (ICACCI), pp 259–265
7. Li Y, Liu Q (2021) A comprehensive review study of cyber-attacks and cyber security; emerging trends and recent developments. Energy Rep 7:8176–8186. ISSN 2352-4847. https://doi.org/10.1016/j.egyr.2021.08.126
8. Bagdatli MEC, Dokuz AS (2021) Vehicle delay estimation at signalized intersections using machine learning algorithms. Transp Res Rec 2675(9):110–126. https://doi.org/10.1177/03611981211036874
9. Farhat YD, Awan MS (2021) A brief survey on ransomware with the perspective of internet security threat reports. In: 2021 9th international symposium on digital forensics and security (ISDFS), pp 1–6. https://doi.org/10.1109/ISDFS52919.2021.9486348

10. Gibson CP, Banik SM (2017) Analyzing the effect of ransomware attacks on different industries. In: 2017 international conference on computational science and computational intelligence (CSCI), pp 121–126. https://doi.org/10.1109/CSCI.2017.20

11. Farion-Melnyk, Rozheliuk V, Slipchenko T, Banakh S, Farion M, Bilan O (2021) Ransomware attacks: risks, protection and prevention measures. In: 2021 11th international conference on advanced computer information technologies (ACIT), pp 473–478. https://doi.org/10.1109/ACIT52158.2021.9548507

12. Homayoun S, Dehghantanha A, Ahmadzadeh M, Hashemi S, Khayami R (2017) Know abnormal, find evil: frequent pattern mining for ransomware threat hunting and intelligence. IEEE Trans Emerg Top Comput 8(2):341–351

13. Jeong D (2020) Artificial intelligence security threat, crime, and forensics: taxonomy and open issues. IEEE Access 8:184560–184574. https://doi.org/10.1109/ACCESS.2020.3029280