



# City Research Online

## City St George's, University of London

**Citation:** Zahoor, U., Rajarajan, M., Pan, Z. & Khan, A. (2022). Zero-day Ransomware Attack Detection using Deep Contractive Autoencoder and Voting based Ensemble Classifier. *Applied Intelligence*, 52(12), pp. 13941-13960. doi: 10.1007/s10489-022-03244-6

This is the accepted version of the paper.

This version of the publication may differ from the final published version. To cite this item please consult the publisher's version.

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/29155/>

**Link to published version:** <https://doi.org/10.1007/s10489-022-03244-6>

**Copyright and Reuse:** Copyright and Moral Rights remain with the author(s) and/or copyright holders. Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge, unless otherwise indicated, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way. For full details of reuse please refer to [City Research Online policy](#).

# Zero-day Ransomware Attack Detection using Deep Contractive Autoencoder and Voting based Ensemble Classifier

*Umme Zahoora<sup>a</sup>, Muttukrishnan Rajarajan<sup>b</sup>, Zahoqing Pan<sup>c</sup>, and Asifullah Khan<sup>\*a,d,e</sup>*

<sup>a</sup>Department of Computer and information Science, Pakistan Institute of Engineering and Applied Sciences, Nilore-45650, Islamabad

<sup>b</sup>School of Mathematics, Computer Science and Engineering, City University of London, London EC1V 0HB, UK

<sup>c</sup>School of Computer and Software, Nanjing University of Information Science and Technology, Nanjing 210044, China

<sup>d</sup>Center for Mathematical Sciences, Pakistan Institute of Engineering and Applied Sciences, Nilore-45650, Islamabad

<sup>e</sup>PIEAS Artificial Intelligence Center (PAIC), Pakistan Institute of Engineering & Applied Sciences, Nilore, Islamabad 45650, Pakistan

Corresponding Author: <sup>\*</sup>Asifullah Khan, [asif@pieas.edu.pk](mailto:asif@pieas.edu.pk)

## ABSTRACT

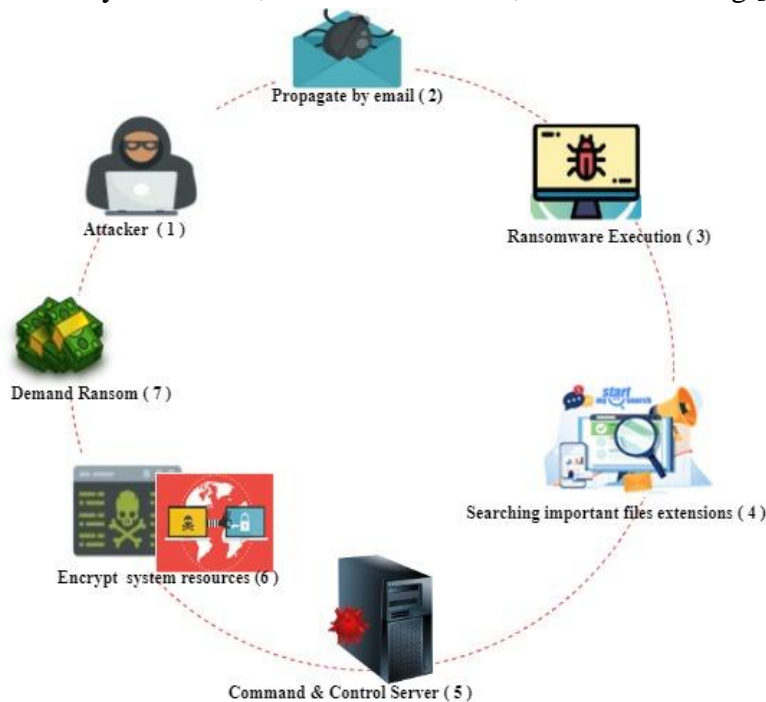
Ransomware attacks are hazardous cyber-attacks that use cryptographic methods to hold victims' data until the ransom is paid. Zero-day ransomware attacks try to exploit new vulnerabilities and are considered a severe threat to existing security solutions and internet resources. In the case of zero-day attacks, training data is not available before the attack takes place. Therefore, we exploit Zero-shot Learning (ZSL) capabilities that can effectively deal with unseen classes compared to the traditional machine learning techniques. ZSL is a two-stage process comprising of: attribute learning (AL) and Inference Stage (IS). In this regard, this work presents a new Deep Contractive Autoencoder based Attribute Learning (DCAE-ZSL) technique as well as an IS method based on Heterogeneous Voting Ensemble (DCAE-ZSL-HVE). In the proposed DCAE-ZSL approach, Contractive Autoencoder (CAE) is employed to extract core features of known and unknown ransomware. The regularization term of CAE helps in penalizing the classifier's sensitivity against the small dissimilarities in the latent space. On the other hand, in case of the IS, four combination rules Global Majority (GM), Local Majority (LM), Cumulative vote-against based Global Majority (CVAGM), Cumulative vote-for based Global Majority (CVFGM) are utilized to find the final prediction. It is empirically shown that in comparison to conventional machine learning techniques, models trained on contractive embedding show reasonable performance against zero-day attacks. Furthermore, it is shown that the exploitation of these core features through the proposed voting based ensemble (DCAE-ZSL-HVE) has demonstrated significant improvement in detecting zero-day attacks (recall=0.95) and reducing False Negative (FN= 6).

---

**Keywords:** Zero-shot Learning, Zero-day Attack, Ransomware, Deep Learning, Autoencoder, Ensemble Classification.

35 **1. Introduction**

36 Ransomware is a malware that possesses its special characteristics in addition to the standard  
37 features of a generic malware [1]. Ransomware, generally follows similar methods to evade,  
38 propagate, and attack its victims as other malwares do. However, it injects its peculiar actions in  
39 the form of processes into target programs, then extract the data, and establishes the connection  
40 with Command and Control server (C&CS). Its main function is to encrypt all the important files  
41 in the target system and demand ransom for recovery. These typical steps of ransomware lifecycle  
42 are described pictorially in Fig. 1. Due to its specific objective, it is considered easy to write and  
43 modify the existing ransomware that can result in an explosive generation of its variants [2].  
44 Wannacry variants (executed in 2017) are said to be responsible for damaging various  
45 organizations that were running an old version of Microsoft Windows. These attacks were  
46 propagated by employing EternalBlue. Some of the other known propagation methods of  
47 ransomware attack include; delivering its payload to the victim using Malicious Emails, bypassing  
48 the typical access control (Bucbi Ransomware [3]), using Exploit Kits (EKs), Injecting redirect  
49 link in JavaScript, Drive-by Download, Waterhole Attacks, and Malvertising [4,5].



50  
51 **Fig 1: Typical steps of ransomware life cycle**

52  
53 The variants of ransomware can be broadly grouped under two categories: "Locker Ransomware"  
54 and "Crypto Ransomware". Locker Ransomware locks the user's system and thus restricts the  
55 access to the system files. Its common variants are: Winlock, DM-4 Locker, CTB Locker, Locky  
56 Ransomware, and Torren Locker. On the other hand, CryptoRansomware, instead of locking the  
57 whole system, encrypts the essential files of the system. Variants of CryptoRansomware include  
58 Pack Crypt, Crypt Locker, Dirty Decrypt, Crypto Wall, and Telsa Crypt. Earlier versions of the

59 Ransomware have used symmetric cryptographic algorithms [6]. However, lately some of its  
60 variants use asymmetric cryptographic methods. Most of the recent Ransomwares are using both  
61 symmetric and asymmetric cryptographic methods that are not easy to break.

62 Existing signature-based detection systems are not able to cope with the increasing number of  
63 unique ransomware variants [7]. Existing intrusion detection systems (IDS) are based on various  
64 analysis method to detect the ransomware. Static analysis based IDS rely on detecting the unique  
65 patterns [8, 9]. Commercially available IDS system adopts static analysis as they are fast and  
66 inexpensive. However, these static IDS normally fails in detecting the zero-day attack and the  
67 polymorphic variations of the attack. In contrast, behavioral based methods focus on the behavioral  
68 profile generated at run time. Behavioral based methods are superior to static analysis in detecting  
69 zero-day attacks and dealing with the polymorphic variations [10]. However, Behavioral based  
70 methods are comparatively slow and cannot detect the metamorphic variations. Anomaly detection  
71 based methods relies on modeling the normal connection behavior and detecting deviations [11,  
72 12]. Such type of anomaly detection methods are better in detecting zero-day attack, but may yield  
73 high false alarm rate.

74 Recently, machine learning based IDS are increasingly used due to their excellent learning  
75 capabilities and their adaptive nature. Researchers exploits various supervised (Decision trees  
76 (C4.5), Support Vector Machines, K-Nearest Neighbor, Naïve Bayes) [13], unsupervised  
77 (Balanced Iterative Reducing and Clustering using Hierarchies (BIRCH)) [14], Deep Learning  
78 (Deep Belief Networks (DBN) [15]), and semi supervised methods [16] to develop intelligent  
79 IDS. Each of these methods has its advantages and disadvantages. However, it is not easy to  
80 achieve the effective IDS for dynamically changing environment with a single classifier and  
81 therefore an ensemble of classifier is more effective. The biggest challenge is to detect zero-day  
82 attacks. Zero-day attacks tries to exploit the new vulnerabilities of the victim's system hence  
83 nothing is known about them in advance. On the other hand, most of the current machine learning  
84 based solutions are dependent upon previous data for detecting future attacks, which is not  
85 available in the case of a zero-day attack.

86 This paper aims to develop a ransomware detection system that can generate an encoding (core  
87 features) based description of the zero-day attacks at run time. Furthermore, it can relate the  
88 derived description with the known attacks for detection purposes. Finally, to increase the  
89 generalization power of the classification system, we have also proposed an ensemble classifier  
90 focusing on the reduction of FN under nominal control of FP. In summary, the proposed technique  
91 is addressing the challenges mentioned above associated with zero-day attacks through the  
92 following contributions:

93 **1.** This work presents a novel zero-shot learning (ZSL) framework to detect the zero-day  
94 ransomware attack. Current approaches generally use an external source of information for

95 attribute learning, which might be a time-consuming task and not viable for collecting data  
96 regarding zero-day vulnerabilities.

- 97
- 98 **2.** The proposed technique presents customized deep contractive autoencoder based attribute  
99 learning (DCAE-ZSL) for zero-day ransomware. For this purpose, an optimum loss function  
100 is learned in an unsupervised manner by optimizing the penalty term for achieving the invariant  
101 representation of the known and unknown ransomware. A lower-dimensional DCAE has been  
102 designed that forces the model to learn only the essential features of the input data. And  
103 compared to traditional ML approaches, it is empirically shown that DCAE based feature  
104 extraction effectively performs well against zero-day Ransomware.
  - 105
  - 106 **3.** An ensemble (DCAE-ZSL-HVE) is trained both on original features and derived attributes to  
107 address semantic loss for suppressing the intra-family variations and increasing the  
108 generalization ability of the classifier on unknown ransomware.
  - 109
  - 110 **4.** The proposed Inference Stage makes cost-sensitive inferences using four simple yet effective  
111 combination rules that provide a considerable compromise between FP and FN.
  - 112
  - 113 **5.** The performance of the proposed technique is compared with both the deep CNN models as  
114 well as shallow learning models.
  - 115

116 The rest of the paper is structured as follows. In Section 2, the related work and its background is  
117 presented. The proposed Attribute Learning (AI) and IS (Inference Stage) phases are presented in  
118 Section 3. The implementation details of the experiments is presented in Section 4. Discussion and  
119 analysis is presented in Section 5. In Section 6, threat to validity of the proposed system is  
120 analyzed. Finally, the conclusion and future work are presented in Section 6.

---

## 121 **2. Related work and background**

122 This section presents a review of recent ransomware detection techniques and their potential in  
123 detecting zero-day attacks. In general existing methods are based on static and dynamic analysis.  
124 Static analysis is performed without executing the malware to extract the structural features.  
125 Andronio et al. [17] presented a static analysis based HelDroid detection system specifically  
126 developed for mobile devices. The model uses only the encryption-based function to detect crypto  
127 and locker ransomware. Mercaldo et al. [18] presented a model-based technique for detecting  
128 mobile-based crypto and locker ransomware. This method analyzes the bytecode of potential files  
129 to inspect only those instructions that are involved in the infection phases. Das et al. [19] presented  
130 a model that performs semantic-based feature extraction by grouping API calls of the same  
131 resource. These API calls of the same resource are further represented as one feature set. The  
132 frequency of these feature sets helps to identify the repetitive actions that differentiates between  
133 benign and malware samples. Alsoghyer et al. [20] proposed an application programming

134 interface (API)-based ransomware detection system (API-RDS) for android platform. The  
135 proposed method identifies the significant API calls from API packages and achieved 97%  
136 accuracy on a self-generated dataset. Although static methods are fast and have a high detection  
137 rate, their main focus is ransomware detection rather than the zero-day attack detection.  
138 Additionally, these methods are generally unable to detect polymorphic variations and packed  
139 families.

140 On the other hand, dynamic analysis is performed by running malware in a safe environment to  
141 extract behavioural features. Kharraz et al. proposed UNEVIL model based on dynamic analysis  
142 [21]. The system's main focus was on learning about the access pattern, system file activities, and  
143 the entropy of I/O data buffer. Song et al. [22] proposed a detection model that performs detection  
144 by monitoring the CPU consumption, memory utilization, files events and I/O usage. Andriono et  
145 al. [17] proposed a dynamic analysis to detect the threatening text. Dynamic analysis based  
146 methods are useful in detecting polymorphic variations. However, they are unable to detect  
147 metamorphic variations. In this regard, few researchers [23–25] have employed a hybrid feature  
148 analysis technique to improve the detection of polymorphic and metamorphic attacks. Alberto  
149 Ferrante et al. [24] proposed a hybrid method to detect mobile ransomware. Firstly, it examines  
150 the potential file using the static method before installation and then observes its runtime behaviour  
151 using dynamic analysis. During static analysis, it computes the frequency of opcodes to detect the  
152 ransomware attack. While during dynamic examination CPU usage, memory consumption, and  
153 network usage are explored to detect malware.

154 Conventional ML-based methods are useful in performing behavioural analysis. The EldeRan  
155 proposed by Sgandurra et al. [26] is based on the files' dynamic behaviour. It uses Mutual  
156 Information Gain to extract the most significant dynamic features and then feeds them to Logistic  
157 Regression Classifier. Hwang et al. [2] proposed, a two-stage mixed ransomware detection model,  
158 based on Markov and Random Forest models. Firstly, their technique builds a Markov model only  
159 on the Windows API call sequence pattern then, builds the Random Forest ML model on the  
160 remaining behavioral features to control FPR and FNR. However, typically the new attacks may  
161 not follow the training distribution to perform covert operations. A different ML-based solution to  
162 zero-day attack detection is an anomaly-based method that trains the model only on normal  
163 activities, and therefore, anything that is not normal is considered as malicious[27]. In [28], Al-  
164 ramy et al. presented a zero-day ransomware detection system using behavioural and data-centric  
165 features. Behavioural features are constructed using n-gram technique on the pre-encryption  
166 generated features. At the same time, data-centric features are generated by grouping the API calls  
167 of the same resources and forming frequency distribution of similar features. After, the feature  
168 generation step, important features are selected using information gain measure. Finally, the  
169 detection module is stacked using two types of classification. Firstly, it performs behavioural  
170 detection using SVM. If the sample is declared as malicious by behavioural detection module, then  
171 it is the final step. Otherwise, the decision is put forward to one class anomaly detection SVM for

172 a final decision. However, the anomaly-based zero-day attack detection system may also generate  
173 a high false alarm rate.

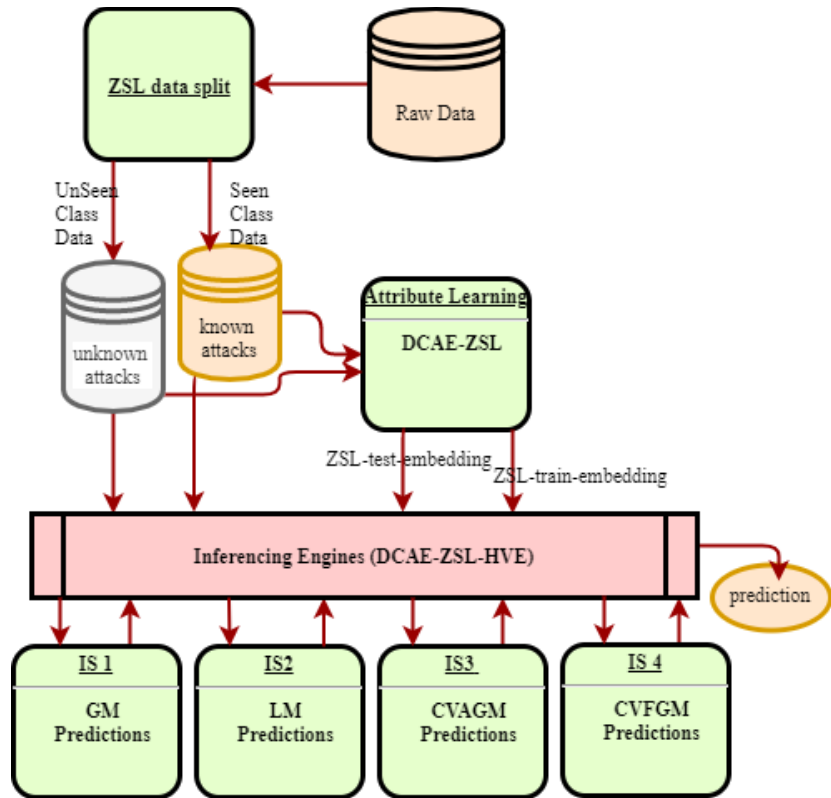
174 Current IDS are mostly unable to cope with modifications in the attack's landscape and are also  
175 highly dependent on the training dataset. On the other hand, an attacker is always searching for  
176 new types of vulnerabilities and exploiting a new weakness in the system, leading to a zero-day  
177 attack. In this regard, we exploit a new paradigm of ML known as *zero-shot learning* (ZSL). ZSL-  
178 based models can detect unseen objects, in the absence of additional knowledge, by relating the  
179 attributes of seen and unseen classes. In literature, interesting techniques have adopted different  
180 methodologies [29] to find the semantic embeddings (core features) between seen and unseen  
181 classes. These embeddings are generally based on attributes that can be derived manually,  
182 discriminatively, using word vectors, mining knowledge from the web, or combining different  
183 kinds of embeddings. ZSL can be described as two-phase process. The first phase is the AL phase,  
184 and the second phase is the IS. In [30], the authors presented a new IS algorithm for Network IDS.  
185 Their main contribution involved an experimental set up for ZSL using of the KDD intrusion  
186 detection dataset. Then, in AL phase applied decision tree for extracting rules. In the IS phase,  
187 they derived the representation depending on point location in Grassmannian manifold, and  
188 explicit distance formula is utilized that finds the shortest distance between the unknown attack  
189 and the known attacks. On the other hand, in [31], XiaoZhang et al. presented regression model  
190 based ZSL method that fits the regression equation for each category. They then, calculated  
191 threshold for all respective categories. In their IS, the test samples' attributes are sequentially  
192 substituted to the equations of all categories. Finally, if the resulted calculation meets the criteria  
193 for all the corresponding thresholds, the attack is considered a known attack; otherwise, it is  
194 considered an unknown attack. Their reported technique can detect unknown intrusion types; such  
195 as Hydra-FTP and HydraSSH types of attacks. However, for a Java-Meterpreter, and Meterpreter  
196 type's other unknown attacks, its capability may not be satisfactory. In [32], Zhang et al., proposed  
197 sparse autoencoders based ZSL method for novel attack detection. It maps known feature space to  
198 semantic space, and try to restore the feature space using reconstruction error constraint. The need  
199 for detecting new ransomware attacks, together with the competence of CAE for change  
200 detection[33], and the competence of ZSL to classify unknown attacks on Zero-shot training data  
201 motivated us to use it for the AL stage. Moreover, high generalization ability of ensemble methods  
202 and effective combination rules encouraged us to use it for the IS.

---

### 203 3. Material and methods

204 This section presents the details of the proposed DCAE-ZSL and DCAE-ZSL-HVE  
205 methodologies. The DCAE-ZSL is the AL method based on finding robust embedding that can  
206 learn the semantic description of zero-day attacks in an unsupervised manner. To learn the context  
207 of a zero-day attack, CAE based feature representation and its effect on ZSL is explored using  
208 zero-day ransomware test data. Finally, the DCAE-ZSL-HVE based is proposed to improve the

209 generalization power and to find a considerable compromise between FP and FN. The abstract  
 210 working of the proposed ZSL model is described in Fig. 2.  
 211



212  
 213 **Fig. 2– Framework of the proposed ZSL architecture**

214 **3.1. Dataset**

215 The ransomware dataset used in the proposed methodologies is accessed from the home page of  
 216 Sgandurra et al. [26] using a given link(<https://www.danielesgandurra.com/>). The data samples  
 217 were retrieved from the VirusShare4 site in the exe files format in February 2016. Later these  
 218 samples were analyzed using Cuckoo Sandbox5 to trace the seven basic features in the runtime  
 219 environment. The dataset consists of 582 ransomware samples and 942 goodware samples.  
 220 Therefore, the available dataset is highly imbalance in nature. The major attributes of the data are  
 221 given in Table 1. Collected samples were further manually categorized into 11 different established  
 222 family’s names. The detail of these families is reported in Table 2. The collected ransomware are  
 223 the most popular variants and mostly are CryptoRansomware. The dataset of goodware is collected  
 224 from trustworthy sources. Goodware application includes browsers, drivers, emulator, and file  
 225 utilities like file search, word office tools, games and various other realistic applications of PC.  
 226 Each sample is analyzed in a sandbox environment for 30 seconds. Although, the authors acquired  
 227 the PCAP traces by connecting VM (Virtual Machines) with a network. However, only host-based  
 228 features were collected.

229 **3.1.1. Features and their visualization**

230 The ransomware uses the API calls to write into the other processes to inject them with its peculiar  
 231 actions or to use it to terminate the other processes. Usually, victim processes are *explorer.exe* or  
 232 *svchost.exe*. On the other hand, the registry keys operation is important because to ensure the  
 233 survivor of ransomware after each reboot. Moreover, it keeps track of every key to get the list of  
 234 mounted devices for exploring more extensions of user files. In literature, ransomware detection  
 235 is often implemented using API call sequences or registry key operations. Other features like DLL  
 236 (Dynamic Link Library) are also important as ransomware often links to the DLL of Visual Basic  
 237 or shell extensions to access the certificates of keys for encryption purpose. Additionally, dropped  
 238 file features are important as ransomware used it for its notes. The important drooped files are html  
 239 or RTF. There are in total of 16382 features collected. These features can be further divided into  
 240 seven broad categories: 1) API represents API invocations, 2) Drop: represents the extensions of  
 241 the dropped files, 3) REG: involves various registration key operations, 4) Files: include operation  
 242 related to files, e.g., create or delete files, 5) FILES\_EXT: is the extensions of all the files that are  
 243 involved in dynamic analysis, 6) DIR: files directory activities, and 7) STR: denotes strings  
 244 embedding. Unlike, the traditional methodologies, the focus of the proposed methods is not on a  
 245 specific action. They are developed by using dynamic features like registry key operations, API  
 246 invocation, files extension, file directory operation, drop files monitoring, files operations and  
 247 embedded strings.

248 **Table 1: Attributes of the dataset**

Parameters	Values
Data type	Binary data
Total malicious samples	582
Total benign samples	942
Total samples	1524
Total features	16382
Missing values	None

249  
 250 **Table 2: Distribution of the different ransomware families**

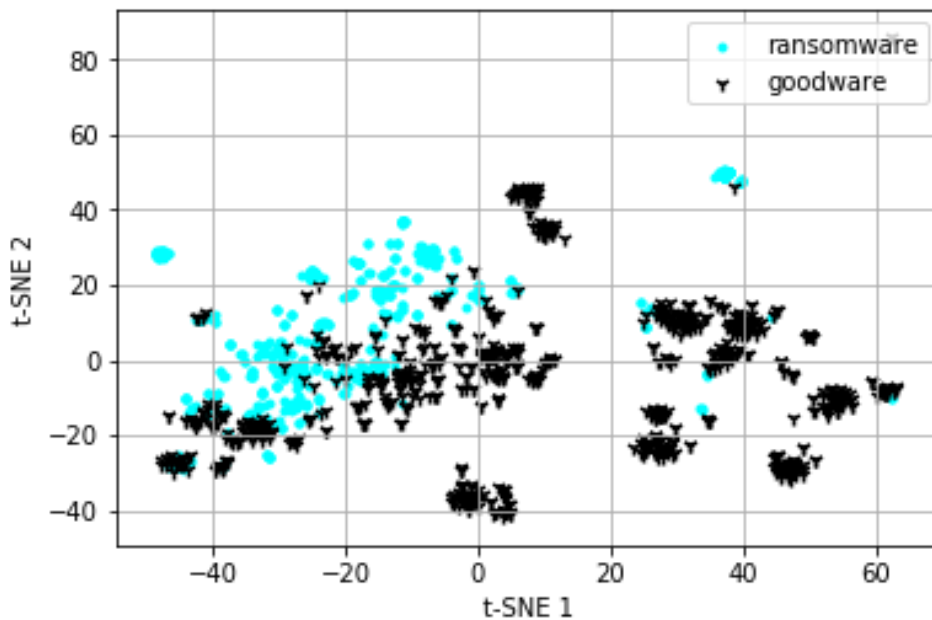
ID	Family Name	Data Distribution	ID	Family Name	Data Distribution
0	Goodware	942	6	Locker	98
1	Critroni	51	7	Matsnu	60
2	CryptLocker	108	8	Pgpcoder	5
3	CryptoWall	47	9	Reveton	91
4	Kollah	26	10	TeslaCrypt	7
5	Kovter	65	11	Trojan-Ransom	35

251 **3.1.2. Data distribution and its visualisation using t-SNE**

252

253 In this work, t-Distributed Stochastic Neighbor Embedding, known as t-SNE [34], is used to show  
254 the data distribution of various classes. t-SNE maps higher dimensional data to two or three-  
255 dimensional space. It is based on SNE optimization and Student’s t-distribution. It finds the  
256 pairwise similarity matrix between the data points. The t-SNE optimization function maintains the  
257 maximum structure of the original mapping. Therefore, in the proposed DCAE-ZSL technique, t-  
258 SNE is used to visualize higher dimensional data into two dimensional lower space, as shown in  
259 Figures 3 and 4. It is useful to visualize the clusters present in malicious and benign data at various  
260 scales. Fig.3. shows the distribution of the ransomware data for the two-class problem. Where the  
261 ‘cyan’ colour is representing the ransomware samples, and the black colour is representing  
262 goodwill samples. It can be observed that the two classes may overlap at some points, thus  
263 indicating the possible similarities between the two categories. There also exist some samples that  
264 are far away from their relative classes due to intra-class variations. Fig.4. shows the distribution  
265 of the multiple families of ransomware more clearly by using different colours.

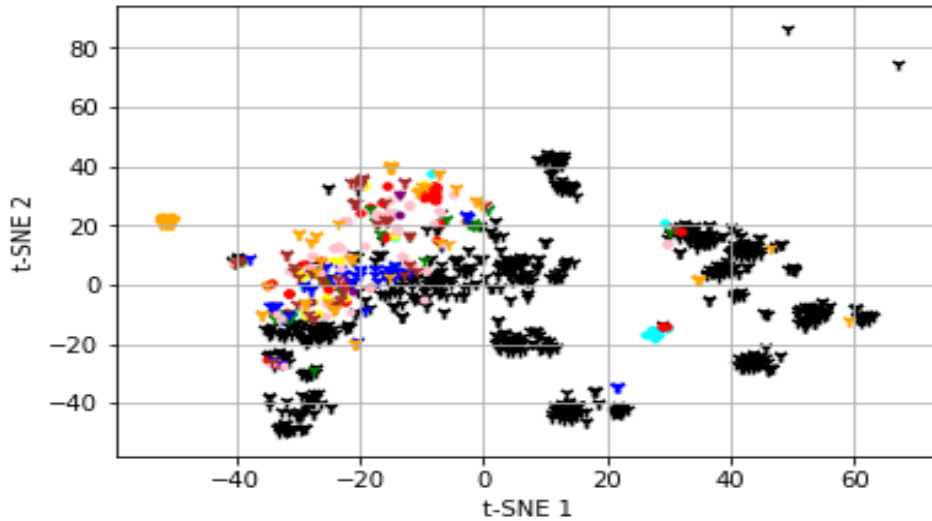
266



267

268 **Fig. 3 Data distribution of ransomware vs goodwill**

269



270

271 **Fig. 4 Data distribution of ransomware families' vs. benign samples. Benign samples are**  
 272 **shown in black color**

273 **3.2. Zero-shot Learning for zero day attack detection**

274 In general, ZSL [29] deals with predicting classes, which were not part of the training. It is mostly  
 275 applicable when the training data does not truly represent all the interested class categories. This  
 276 situation is most viable for IDS due to the growing numbers of zero-day attacks. Correspondingly,  
 277 in ZSL, classes of training and test instances are different from each other. Classes of training  
 278 instances are known as seen classes, and test instances are known as unseen classes. ZSL can be  
 279 described as a two-stage process: AL and IS. Attributes are learned using derived knowledge from  
 280 some labelled training data or finding some intermediate information to relate the seen and unseen  
 281 classes. In the IS, the derived knowledge is used to detect the unseen classes. Therein literature,  
 282 different AL algorithms are presented in the context of IDS, e.g. Attribute Learning for Network  
 283 Intrusion Detection (ALNID)[31], Graph Embeddings[32], Deep Attribute Prediction  
 284 (DeepAP)[33], and Grassmannian [24]. ZSL can be classified as Inductive ZSL and Transductive  
 285 ZSL subject to the available information. Inductive ZSL learned attributes using only seen class  
 286 information, while Transductive ZSL uses both labelled unseen data and unlabeled unseen data.  
 287 Based on a test set, ZSL is classified as conventional ZSL (CZSL) or Generalized ZSL (GZSL).  
 288 CZSL evaluates its model on unseen classes, whereas, GZSL uses both seen and unseen classes in  
 289 evaluation. The proposed DCAE-ZSL method is the transductive approach. It generalized the  
 290 model on unknown attacks by utilizing both seen and unseen classes in an unsupervised manner  
 291 and evaluated using the CZSL method.

292 **3.3. ZSL-Data Split**

293 We firstly performed data partitioning for training and evaluation of the ZSL in ransomware  
 294 detection related tasks to achieve the desired goals. For this purpose, we split the original data 'D'

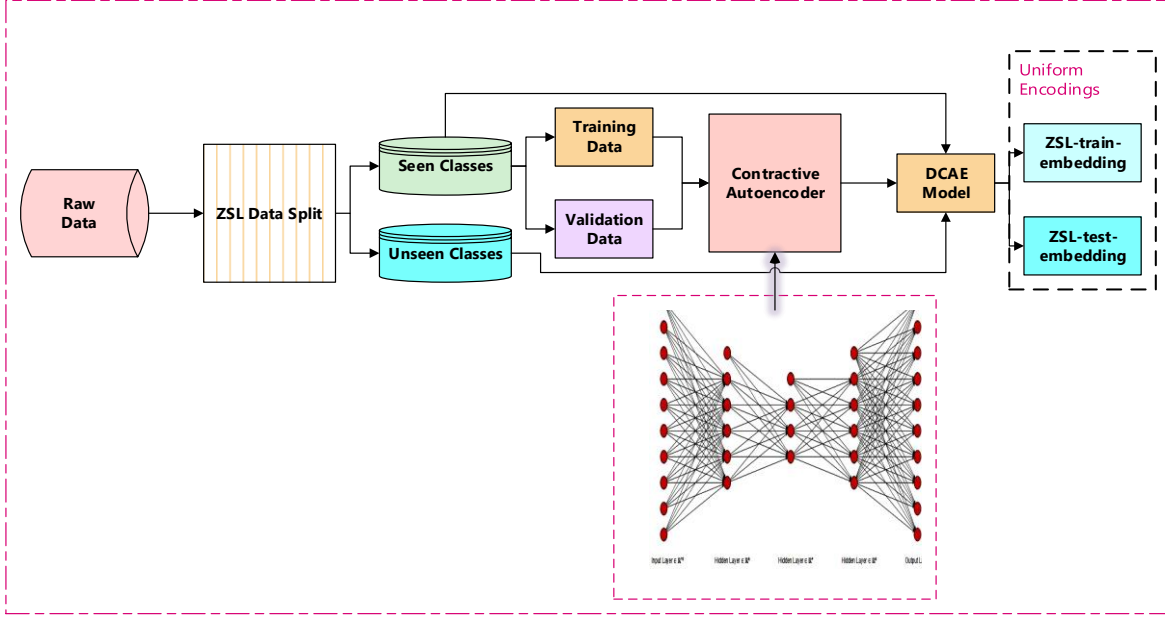
295 into two disjoint sets of seen and unseen classes, named as ZSL-train data =  $\{X_i^r, Y_i^r\}$  and, ZSL-  
 296 test data =  $\{X_i^s, Y_i^s\}$ , respectively. Where ‘ $X_i$ ’ is representing the attributes of the samples, and  
 297 ‘ $Y_i$ ’ is representing their respective labels. The proposed system presents the binary classification  
 298 of ransomware vs goodware( $G_i^n$ ). Models are trained on  $Y^r$  members and tested on  $Y^s$  members.  
 299 Respective members of  $Y^r$  and  $Y^s$  are described as sets in (Eq.1 and Eq.2 respectively), where  
 300  $Y^r \cap Y^s = \emptyset$ . In the context of ransomware detection seen classes are known attacks on which  
 301 models are trained. However, unseen classes are unknown zero day ransomware on which models  
 302 are tested.

$$303 \quad Y^r = \{G_i^k, \text{Critroni, CryptLocker, CryptoWall, Kollah, Kovter, Locker, Matsnu}\} \quad (1)$$

$$304 \quad Y^s = \{G_{k+1}^n, \text{Pgpcoder, Reveton, TeslaCrypt, Trojan-Ransom}\} \quad (2)$$

### 305 **3.4. The proposed DCAE-ZSL attribute learning method**

306 This module's objective is to generate the class independent description of seen and unseen (zero-  
 307 day) attacks. Autoencoder can learn useful latent representations without class information. In  
 308 essence, the learned representation should represent the core features by suppressing unnecessary  
 309 variations. CAE is used to suppress unnecessary variations. Further, Deep Undercomplete  
 310 Autoencoder is used to extract the most useful core representation of the data. To build DCAE-  
 311 ZSL, we have trained ten hidden layers of CAE on ZSL-train data. Where the first five layers are  
 312 encoding layers, and the remaining five are the decoding layers. However, the DCAE-ZSL  
 313 technique selected the robust features from the fifth encoding layer comprised of 100 neurons.  
 314 Overall, the proposed DCAE-ZSL methodology is schematically explained in Fig.5.  
 315 Implementation details of designed topology are described in section 4.1. Whereas, Table 3 is  
 316 illustrating the designed topology of the CAE.



317

318 **Fig. 5– Overview of the proposed Attribute Learning (AL) architecture (DCAE-ZSL)**

319 **Table 3: Number of neurons in encoding and decoding layers of Autoencoder**

	E1	E2	E3	E4	E5	D6	D7	D8	D9	D10
<b>Input Unit</b>	16382	4000	2000	1200	600	100	600	1200	2000	4000
<b>Output Unit</b>	4000	2000	1200	600	100	600	1200	2000	4000	16382
<b>Activation Function</b>	Relu	Relu	Relu	Relu	Relu	Relu	Relu	Relu	Relu	Sigmoid

320 **3.4.1. Autoencoder**

321 Autoencoder is an unsupervised neural network that backpropagates by setting its input as its target  
 322 value [35]. A simple Autoencoder is composed of an input layer, a hidden layer, and an output  
 323 layer. The Autoencoder's objective is to learn useful hidden representations with minimum  
 324 reconstruction loss.

325 
$$L = \|x - x'\|^2 \tag{3}$$

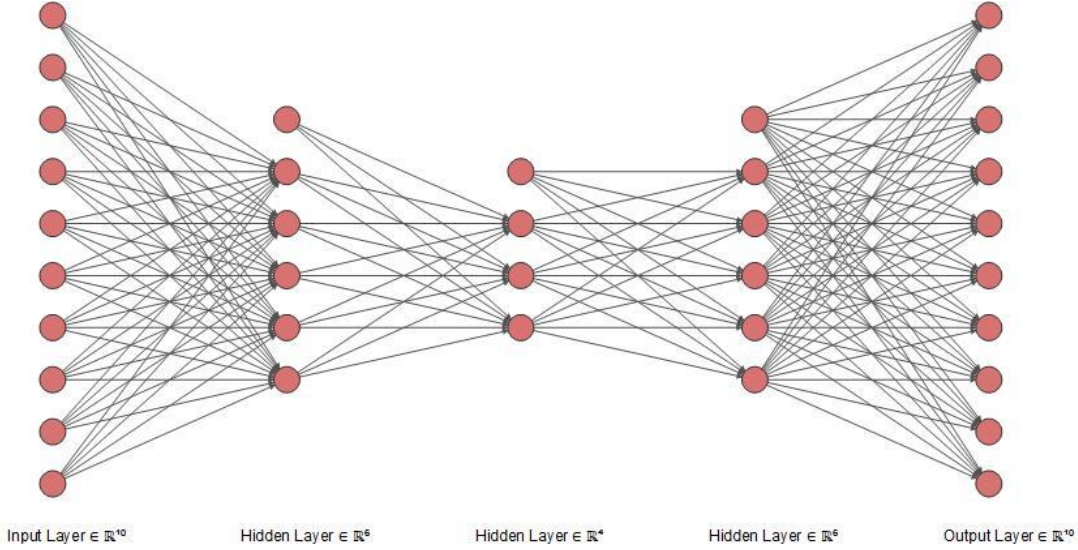
326 Autoencoder based machine learning algorithm is a two-step process that involves pre-training  
 327 and fine-tuning [36]. Pre-training is unsupervised learning that backpropagates error in a greedy  
 328 layer-wise manner to reconstruct its input with minimum loss. Pre-training consists of the encoding  
 329 and decoding layers. In encoding, each layer learns important features from previous information;  
 330 subsequently, these essential features are the original input's encoded form. The mathematical  
 331 description of the encoding process thus becomes:

332 
$$Encoding = x' = \partial(W_e x + b_e) \tag{4}$$

333 While in decoding layers, each encoded layer is decoded to reconstruct the original features. The  
 334 process of decoding is mathematically described in Eq. (5):

$$335 \quad \text{Decoding} = x = \hat{\sigma}(W_d x' + b_d) \quad (5)$$

336 where,  $x$  and  $x'$  are representing the original and encoded signals, respectively,  $W_e$  is the weight  
 337 matrix of encoding layers and  $W_d$  is the weight matrix for the decoding layers,  $b_e$  and  $b_d$  are their  
 338 respective biases, and  $\hat{\sigma}$  is the activation function to add the nonlinearity.



339

340 **Fig. 6 Undercomplete Autoencoder**

### 341 3.4.2. Contractive Autoencoder (CAE)

342 CAE is an unsupervised regularized autoencoders variant[37]. The CAE proposed by Rafie et al.  
 343 is the extension of the Denoising (DE) autoencoder having the same motivation to produce a robust  
 344 representation. It allows small perturbation around its training data using an additional penalization  
 345 term ( $\|J_h(x)\|_F^2$ ) in its generic loss function. This penalization term penalizes the large derivative  
 346 using a lambda parameter that controls the extent of change in input, w.r.t the learned  
 347 representation. This penalization term forces the model to learn the uniform representation of the  
 348 concept (as described in Eqs. (6) and (7)).

$$349 \quad L = \|x - x'\|^2 + \lambda \|J_h(x)\|_F^2 \quad (6)$$

$$350 \quad \|J_h(x)\|_F^2 = \sum_{ij} \left( \frac{\partial h_j(x)}{\partial x_i} \right)^2 \quad (7)$$

351 CAE can learn the non-linear manifold. The common variations present in the data correspond to  
 352 the manifold's local dimensions. Whereas, the variations that are of rare type will correspond to  
 353 the orthogonal dimension. The penalization term ensures invariant feature along all dimensions.  
 354 At the same time, the reconstruction error term offers to reconstruct the input faithfully. In case

355 of directions having strong contractive pressure are the ones, where the input density is sparse. In  
356 comparison, the directions with weak contractive pressure raise the input density. Hence, training  
357 set direction can resist this contractive pressure.

358 Alternatively, Contractive mapping can be achieved by weight decaying in the linear case. For  
359 non-linearity, it can be achieved by motivating the hidden units to their saturated regime. Sparse  
360 autoencoders achieve it by keeping most of its components close to zero. Thus, maintains tiny  
361 derivatives in its Jacobian term. Similarly, DE autoencoders indirectly incorporate the robustness  
362 in its reconstructions phase. They achieved scholastically by generating the corrupted input to  
363 attain identity function, while reconstructing the clean version. Whereas, CAE's are the ones that  
364 explicitly encourage robustness in its encoding phase, which is more vital as it penalizes the  
365 magnitude of the first derivatives at training data. This property makes it as a suitable choice for  
366 feature extraction purposes as compared to others.

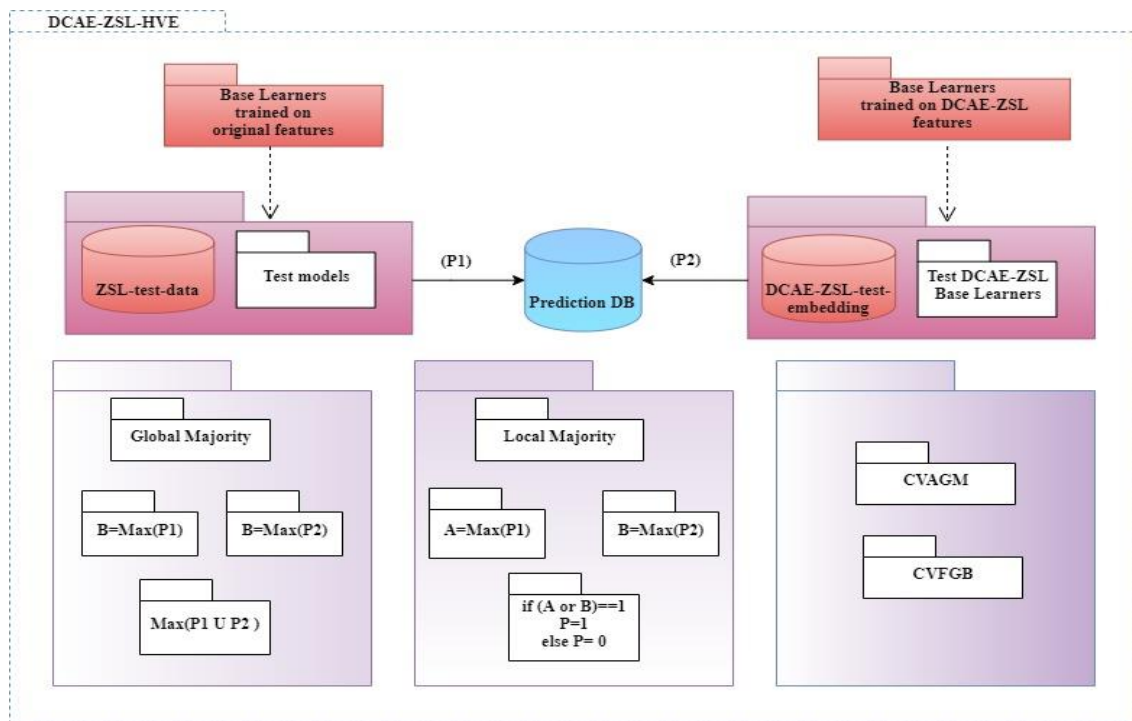
### 367 **3.5. Inference Stage (IS)**

368 In this stage, classes are inferred from the extracted attributes. Therein literature has three well-  
369 known IS methods [38]: Probabilistic frameworks [39], energy function [38, 40], and the K-  
370 Nearest Neighbor (K-NN). The probabilistic framework's most general form combines the derived  
371 attributes with original attributes to find the target class. However, it has two variants: (i) Directed  
372 Attribute Prediction (DAP) and Indirected Attribute Prediction (IAP). DAP method builds the  
373 learning model for each attribute of AL stage. These learning models are then used in IS to predict  
374 new classes using attribute signature. IAP creates the learning model for each training class. At  
375 the evaluation stage, the prediction from each training class tempt labelling of the attribute layer,  
376 which helps to infer the test class labels[41]. The K-Nearest Neighbour (K-NN) based methods  
377 rely on distance measures to find the closest match of input instance to the attribute instance  
378 derived in AL stage. The proposed IS method learns various models e.g.: Random Forest (RF),  
379 Logistic Regression (LR), Support Vector Machine (SVM), Gaussian Naïve Bayes (GNB) on  
380 attributes learned in DCAE-ZSL stage and on original attributes. At test time, the two feature  
381 representation schemes' inferences are combined using GM, LM, CVAGM, and CVFGM on  
382 unseen classes of attack. The proposed (DCAE-ZSL-HVE) ensemble scheme with its combination  
383 rules is described below.

#### 384 **3.5.1. Proposed ensemble (DCAE-ZSL-HVE) based Inference Stage**

385 In this paper, we applied different voting-based ensemble approaches to infer unknown  
386 ransomware. Firstly, the Global Majority (GM) based method selects the majority votes from all  
387 the decision models considering them as individual decision spaces. Algorithm 1 describe the GM  
388 procedure. Besides, two decision spaces are distinctly defined based on features representation  
389 scheme in the Local Majority case (Algorithm 2 show the LM procedure). With this voting-based  
390 mechanism, a winner class label is selected based on the local majority of any single decision  
391 space. It is experimentally observed that multiple base learners can achieve greater

392 generalizability. However, not all the base learners can yield effective classification performance  
 393 in a practical pattern recognition task regarding the diversity or pseudo-independent nature of base  
 394 learners. Some base learners may have adverse effects on the ensemble learner. Therefore, we  
 395 carried out a model selection by finding the cumulative vote-against value for each model.  
 396 Top 'n' models with the minimum cumulative vote against them are selected to form a single  
 397 decision space in case of global majority voting (CVAGM). Similarly, in the case of the proposed  
 398 CVFGM cumulative vote for value is calculated for each model. Then select the top 'n' models  
 399 with maximum cumulative vote-for value. Detailed procedure of the combination schemes  
 400 CVFGM and CVAGM is described in Algorithm 3. The aim was to reduce the FN by increasing  
 401 the generalization of ensemble and giving more weightage to FN. Overview of the proposed  
 402 (DCAE-ZSL-HVE) architecture is described in Fig.7.



403

404 **Fig. 7– Overview of the proposed (DCAE-ZSL-HVE) Inference Stage (IS) architecture**

---

**Algorithm 1 : Local Majority based Combination scheme**

---

**Input:**Local-DB1: $m \times n$  dimensional, prediction matrix of base estimators trained on original features

Local-DB2: $m \times n$  dimensional, prediction matrix of base estimators trained on contractive features

**Output :** Final-Prediction

1.  $n \leftarrow$  no-of-estimators;
  2.  $m \leftarrow$  test-sample-size;
  3. Final-prediction  $\leftarrow []$ ;
  4. **for**  $r$  in range( $m$ ):
  5.      $A = \text{Max}(\text{Local-DB1}[r, :])$ ;
  6.      $B = \text{Max}(\text{Local-DB2}[r, :])$ ;
  7.     **if** ( $A == 1$  or  $B == 1$ );
  8.         Final-prediction[ $r$ ]  $\leftarrow 1$ ;
  9.     **else**
  10.         Final-prediction[ $r$ ]  $\leftarrow 0$ ;
  11.     **end for**
  12. **return** Final-Prediction;
- 

405

---

**Algorithm 2 : Global Majority based Combination scheme**

---

**Input:**Local-DB1: $m \times n$  dimensional, prediction matrix of base estimators trained on original features

Local-DB2: $m \times n$  dimensional, prediction matrix of base estimators trained on contractive features

**Output :** Final-Prediction

1.  $n \leftarrow$  no-of-estimators;
  2.  $m \leftarrow$  test-sample-size;
  3. Final-prediction  $\leftarrow []$ ;
  4. Global-DB = Local-DB1  $\cup$  Local-DB2;
  5. **for**  $r$  in range( $m$ ):
  6.     Final-prediction[ $r, :$ ]  $\leftarrow \text{Max}[\text{Global-DB}[r, :]]$ ;
  7.     **end for**
  8. **return** Final-Prediction;
- 

406

**Algorithm 3 : Commutative Vote For & Cumulative Vote Against based Global Majority Combination Schemes**

**Input:** *Local-DB1*:  $m \times n$  dimensional, prediction matrix of base estimators trained on original features // The prediction matrix is in the form of zero and one ,where '1' is representing ransomware and '0' is representing goodware.

*Local-DB2*:  $m \times n$  dimensional, prediction matrix of base estimators trained on contractive features

**Output :** Final-Prediction

```
1.  $n \leftarrow$  no-of-estimators;
2.  $m \leftarrow$  test-sample-size;
3.  $k \leftarrow$  No-of-top-classifiers;
4. Global-DB =Local-DB1 U Local-DB2;
5. for  $r$  in range( $m$ ): // Traverse Global-DB
6.   cumulative_vote_for  $\leftarrow$  [];
7.   cumulative_vote_against  $\leftarrow$  [];
8.   for  $i$  in range( $n$ ):
9.     No_of_Vote_for [ $r$ , $i$ ]  $\leftarrow$  Count no of predictions similar to the  $i^{\text{th}}$  base estimator;
10.    No_of_Vote_against [ $r$ , $i$ ]  $\leftarrow$  Count predictions dissimilar to the  $i^{\text{th}}$  base estimator;
11.   end for
12. end for
13. for  $j$  in range( $m$ ):
14.   for  $t$  in range( $n$ ):
15.     cumulative_vote_for[ $t$ ] = sum( No_of_Vote_for [ $j$ ,:]);
16.     cumulative_vote_against[ $t$ ] = sum( No_of_Vote_against [ $j$ ,:]);
17.   end for
18. end for
19. Select top  $k$  base learners from cumulative_vote_for matrix;
20. Apply majority voting on selected base learners;
21. return Final-Prediction; // predictions by CVFGM;
22. Select top  $k$  base learners from cumulative_vote_against matrix;
23. Apply majority voting on selected base learners;
24. return Final-Prediction; // predictions by CVAGM;
```

407

408 **3.6.Performance comparison with state of the art Deep Models**

409 Deep learning (DL) techniques are becoming famous for enhancing the performance of network  
410 intrusion detection systems (NIDS)[42]. However, DL based detection models are highly dependent  
411 upon large amounts of labelled training data, often unavailable. In this comparison, we investigate the  
412 potential of transfer learning (TL) in detecting zero-day ransomware. TL is an approach that enables  
413 the transferring of learned features from the source domain to the target domain, especially in cases  
414 where the target domain samples are less in number. We compare the performance of the proposed  
415 model with several deep models trained using TL, e.g. ResNet50[43], GoogleNet (Inception-V1)[43],  
416 and Inception-V3[44]. We used the pre-trained model trained on the imageNet dataset and then fine-  
417 tuned the model on ZSL-train data. The implementation details of these architectures are provided in

418 the section 4.2. These TL based deep CNN architectures are customized to make them applicable for  
419 ransomware image dataset by adjusting new input layer as per the dimensions of the targeted  
420 ransomware image data (224X224X3). Similarly, the last fully connected layer of the standard  
421 architecture is replaced with two neurons to classify the goodware and ransomware samples using  
422 softmax function at the last layer

423 On the other hand, the remaining convolutional blocks are kept unchanged. These models' weight space  
424 is optimized using the backpropagation algorithm to minimize the cross-entropy-based loss function.  
425 We show how effective are the TL based detection models in detecting zero-day attacks, when there is  
426 no information available, compared to the proposed and shallow learning methods. GoogleNet is a  
427 Deep CNN architecture that employs inception block to transforms the image representation at multiple  
428 scales using multi-resolution filters. However, Inception-V3 replaces these multi-resolution filters with  
429 asymmetric filters to make them computationally efficient. ResNet improves the optimization strategy  
430 for fast convergence using skip connections.

---

## 431 **4. Implementation details**

432 All the experiments were carried out using Microsoft Window7 professional, 16.0 GB memory,  
433 X64-based PC, 64-bit operating system. The core coding modules are implemented using python  
434 version 3.6. However, the base learners SVM, RF, GNB and LR were developed using Scikit  
435 library. Proposed CAE topology is designed using Keras library. In addition, pillow and t-SNE  
436 libraries were also used.

### 437 **4.1. Parameter settings of baseline models and Deep Contractive Autoencoder**

438 Parameter setting involves the parameters of the DCAE and different models, including RF,  
439 GNB, SVM, and LRC. To extract DCAE bottleneck features, CAE is trained on the ZSL-train set.  
440 To optimize the parameters, 15% of the total data is reserved as validation data. Table 4 shows the  
441 values of the parameters, which are set during the training phase of DCAE. CAE is trained using  
442 a contractive-loss function that takes  $\lambda = 0.000001$ . Different conventional learning model's  
443 optimization is carried out using 5-cross validation on precision and recall due to the data's  
444 imbalanced nature.

445  
446

**Table 4: Parameter setting of Autoencoder**

<b>Parameters</b>	<b>Values</b>
Total encoded layers	4 + 1
Total decoded layers	5
Batch size	2
Encoding dimension	100
Total layers	10
Lambda	0.000001

## 4.2. Parameter settings of TL based CNN models

TL-based CNN models' parameter setting involves its hyperparameters settings like Epochs: 100, learning rate: 0.0001, batch size: 5, and momentum: 0.95. Furthermore, the train, test and validation data are kept the same for comparison purpose both for the baseline models and the proposed technique. TL-based Deep CNNs models' training is optimized using SGD optimizer to minimize the cross-entropy loss.

## 4.3 Performance evaluation

To evaluate the proposed DCAE-ZSL technique's effectiveness, the different performance measures used in various experiments are recall, accuracy, precision-recall curve (PR Curve), and AUC-PR. PR Curve is a plot between precision and recall by varying threshold. The mathematical description of the used performance measures is given below in Eqs. (8-10).

$$DetectionRate = Precision = \frac{TP}{TP + FP} \quad (8)$$

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \quad (9)$$

$$Recall = \frac{TP}{TP + FN} \quad (10)$$

The positive (P) and negative (N) represent two classes, in the current problem, the positive class is the ransomware, and the negative is goodware. The objective is to increase TP value and decrease FN value, as the malicious class detection is more important in this scenario.

---

## 5. Results and discussion

### 5.1. Effectiveness of the proposed DCAE-ZSL attribute learning method

The objective of this experiment is to show that the robust features extracted through the bottleneck layer of Undercomplete CAE can work well against zero-day attacks. Table 5 shows the result of zero-day attack detection using a full set of original features. In addition, the results of zero-day attack detection using robust and reduced features extracted through the proposed DCAE-ZSL technique. It illustrates the obtained results in terms of the recall and accuracy. DCAE-ZSL-RF has shown significant improvement in recall (0.85) and accuracy (92.8) as compared to recall (0.79) and accuracy (81.0) of the baseline RF model that was trained on the original features. Similarly, DCAE-ZSL-GNB and DCAE-ZSL-SVM have shown considerable improvement in terms of the recall and accuracy, as compared to the baseline GNB and SVM model. DCAE-ZSL-LR performance has also shown improvement in terms of the recall and is more sensitive in detecting the positive class than the baseline LR model. Overall it can be observed that features extracted through the proposed technique can enhance the ZSL capabilities of conventional learning algorithms.

479 **Table 5:** Different classifiers’ result on original and transformed feature space (proposed DCAE-  
 480 ZSL) on test set.

Attributes Type	Learning Models	TN	FP	FN	TP	Recall	Accuracy(%)
<b>Original</b>	<b>RF</b>	112	21	28	106	0.79	81.0
<b>Proposed DCAE-ZSL</b>		133	0	20	115	<b>0.85</b>	<b>92.8</b>
<b>Original</b>	<b>GNB</b>	105	28	47	87	0.64	71.9
<b>Proposed DCAE-ZSL</b>		84	49	11	123	<b>0.91</b>	<b>77.5</b>
<b>Original</b>	<b>SVM</b>	130	3	28	106	0.79	88.4
<b>Proposed DCAE-ZSL</b>		120	13	13	121	<b>0.90</b>	<b>90.3</b>
<b>Original</b>	<b>LR</b>	132	1	25	109	0.81	90.4
<b>Proposed DCAE-ZSL</b>		121	12	14	120	<b>0.90</b>	<b>90.3</b>

481

482 **5.2 Effectiveness of the proposed ensemble (DCAE-ZSL-HVE) based Inference Stage**

483 We present the extensive experimental results for different ensemble configurations using different  
 484 learning models (RF, DCAE-ZSL-RF, SVM, DCAE-ZSL-SVM, LR, DCAE-ZSL-LR, DCAE-  
 485 ZSL-GNB and GNB) trained distinctly on two feature representations in Table 6. The bold values  
 486 indicate the best values for the each method used. Table 6 shows the performance of diverse voting  
 487 approaches that define the different combination rules on the dynamic behavior dataset. The  
 488 highest recall performance is obtained by using the local majority (i-e 0.95). This combination rule  
 489 and ensemble scheme provide a 92.8% accuracy rate. Besides, the global majority based  
 490 combination yields the second greatest results in terms of recall (i-e. 0.91). So, we conclude that  
 491 the local majority can gain a recall measure for ransomware detection. The performance of the  
 492 four different proposed ensemble is evaluated using recall. FP, FN, accuracy and F1 show that  
 493 each inferencing method brings improvements than the individual learners.

494

495

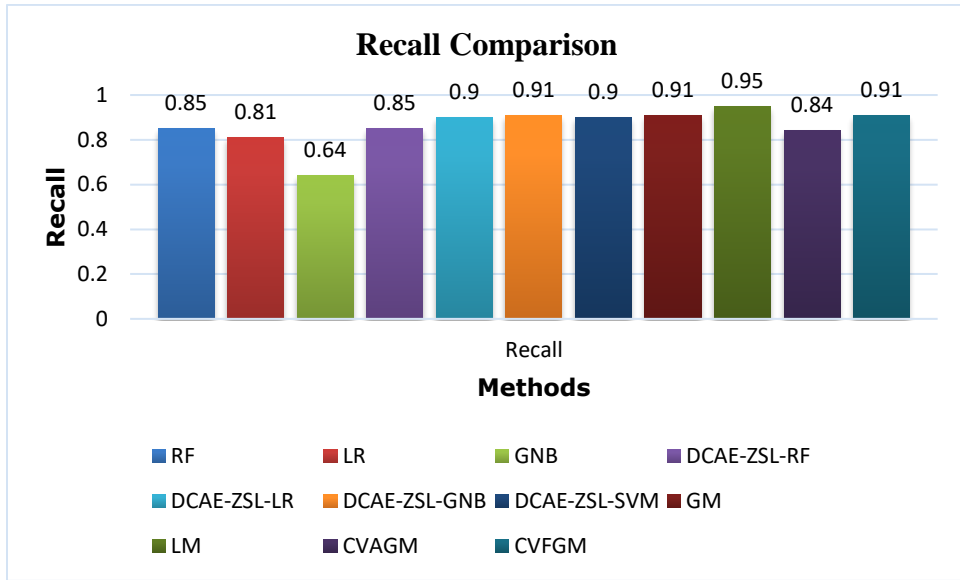
496

497 **Table 6:** Comparison of DCAE-ZSL-HVE with four different proposed combination rules in terms  
 498 of Recall, Accuracy, and F1 on test data

IS Methods	Combination Schemes	TN	FP	FN	TP	Recall	Accuracy (%)	F1
Proposed Method I	GM	129	4	11	123	0.91	<b>94.3</b>	<b>0.94</b>
Proposed Method II	LM	120	13	<b>6</b>	<b>128</b>	<b>0.95</b>	92.8	0.93
Proposed Method III	CVAGM	<b>131</b>	<b>2</b>	21	113	0.84	91.0	0.91
Proposed Method IV	CVFGM	120	13	11	123	0.91	91.0	0.91

499  
 500 **5.3 Performance analysis using recall**

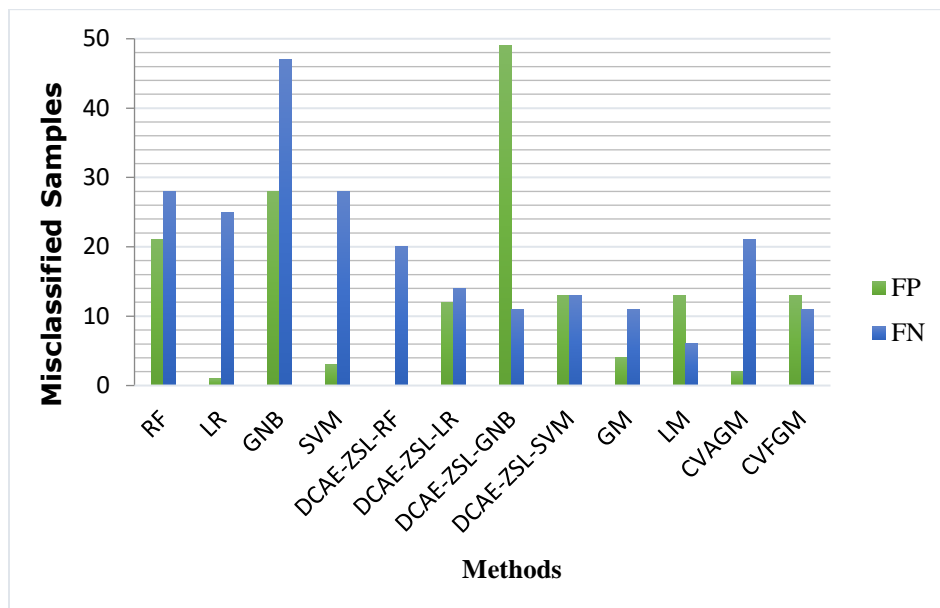
501 Fig.8. shows the recall measure comparison of, baseline models, proposed DCAE-ZSL based  
 502 models and the proposed ensemble methods. From the figure, it can be seen that the proposed  
 503 DCAE-ZSL transformation brings improvement in terms of recall for all learning models.  
 504 Moreover, the proposed ensemble scheme further improves the recall rate. The results of all  
 505 combination rules except Method III are better than the base learners trained on the two  
 506 representation schemes. However, the best results (0.95) in terms of recall are obtained by using  
 507 the local majority as a combination rule.



508  
 509 **Fig. 8– Performance comparison in terms of recall**

#### 510 5.4 Performance analysis using FP and FN

511 In this section, we further evaluate the FP and FN error measures. Fig. 9 shows the True Positive  
512 (TP), and the FP values, of the classifiers trained on original features, proposed DCAE-ZSL  
513 transformation and trained on proposed ensemble scheme with four different combination rules.  
514 The figure shows that the lowest FP is achieved by the proposed DCAE-ZSL-RF model (FP=0).  
515 As this research work aims to decrease the FN value, it can be observed that it is decreasing with  
516 each proposed module. The lowest FN is achieved using the proposed ensemble scheme with the  
517 local majority combination rule (i.e., FN =6). FN values of all the classifiers that are trained on the  
518 proposed transformation are better than the baseline models which are trained on original features.  
519 When applied the proposed ensemble, it further decreasing the FN values.



520

521 **Fig. 9– Performance comparison in terms of FP and FN values**

#### 522 5.5 Detection rate analysis at multiple threshold levels

523 The ransomware attack is one of the most dangerous cyber-attacks due to its common effects on  
524 internet resources. Therefore, false alarms pose less cost as compared to a FN. Hence, we  
525 optimized the recall by compromising precision, as this compromise is desirable. Fig.10. (a) and  
526 (b) are showing the PR curves of baseline and proposed DCAE-ZSL techniques. It can be observed  
527 that when the recall threshold is less than 0.8, the baseline is performing better. However, as the  
528 baseline approaches 0.8, the proposed techniques outperform both in terms of precision and recall.  
529 Even below 0.8 recall threshold values, results are comparable.

530

531

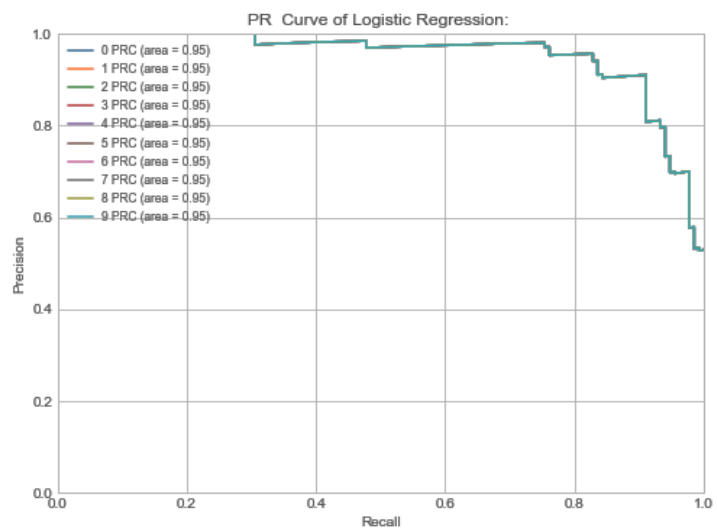
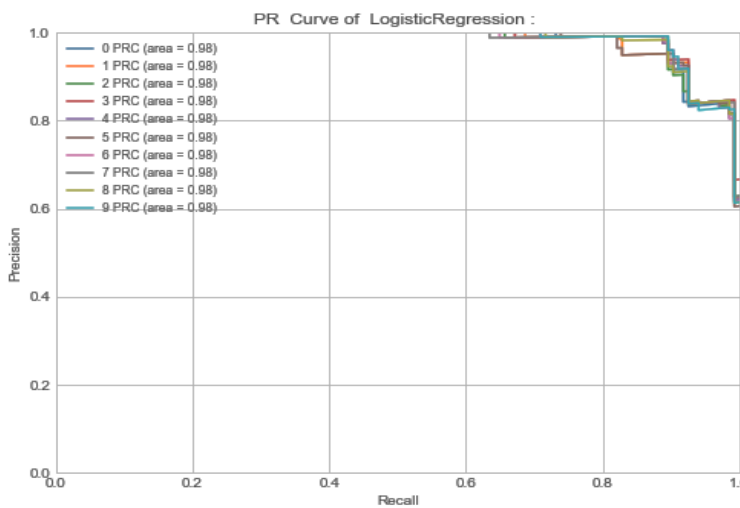
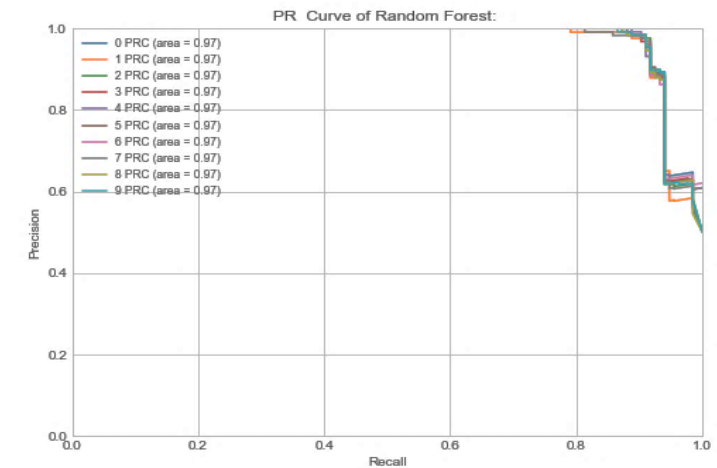
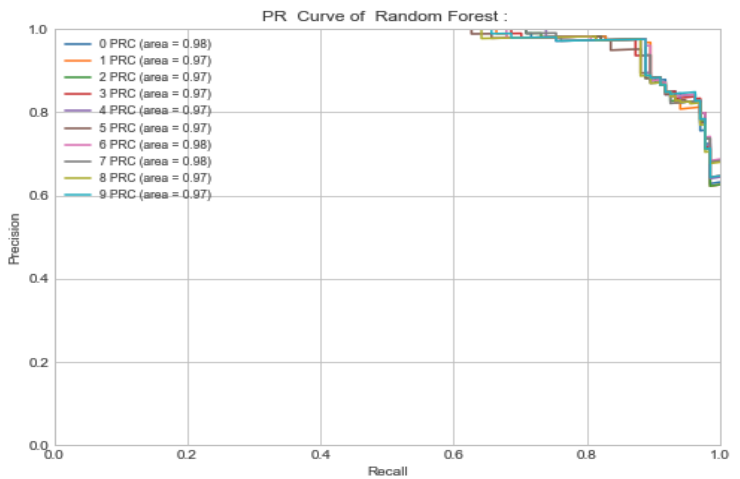
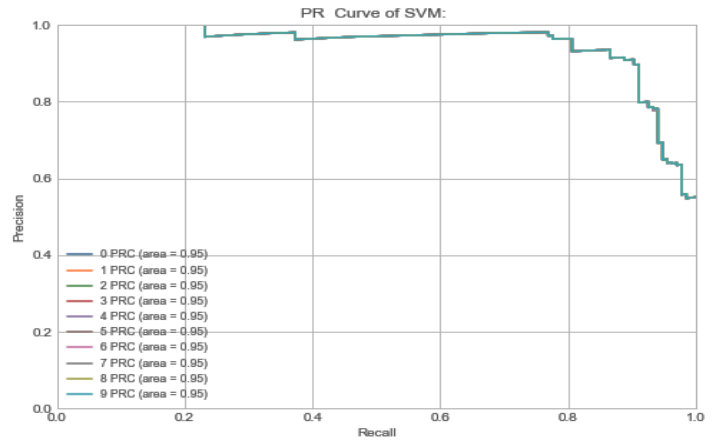
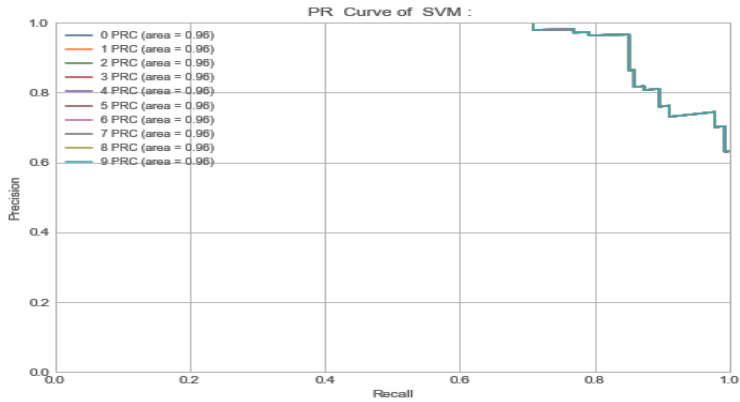
532

533

534  
535  
536

### Original Attributes(Full Set)

### Proposed DCAE-ZSL Attributes(100 Features)



537

Fig. 10– PR curve of baseline vs proposed technique

538

539 **5.6 Performance Comparison with state of the art techniques**

540

541 Table 7 shows the comparison between the proposed method's performances and fine-tuned CNN  
 542 models on ZSL-test data. The performance of the proposed DCAE-ZSL-HVE is compared with  
 543 three classification models (GoogLeNet, Inception- V3, ResNet-50) on the same data using the  
 544 same train test distribution. These classification models are now in trend to solve a complex  
 545 problem like intrusion detection[45]. The Performance analysis of TP, FN and recall values  
 546 suggests that proposed DCAE-ZSL-HVE models learn the ransomware specific feature better than  
 547 GoogLeNet, Inception- V3, ResNet-50. In comparison, the fine-tuned CNN models show slightly  
 548 better performance in learning goodware specific features but at the cost of a high false alarm rate  
 549 (minimum 23 by Inception- V3). However, the proposed DCAE-ZSL-HVE generates only 13 false  
 550 alarm. This may be because GoogLeNet, Inception- V3, ResNet-50 are data-hungry classifiers and  
 551 are biased towards majority class samples.

552

553 Fig.11. depicts the Bar Chart graph comparison among best performers of baseline models (i-e  
 554 LR), Deep CNN models (i-e ResNet50), the proposed IS (i-e LM) in term of error measures (FP  
 555 and FN), recall value and number of misclassified samples (MCS). The number of MCS is  
 556 calculated by summing the error measures FP and FN. It can be observed that the values of all  
 557 metrics using the proposed ensemble-based IS (LM) is better than all the other reported results,  
 558 except FP for baseline LR. However, LR yields high MCS than the proposed IS (LM). This  
 559 indicates that the proposed technique achieves a considerable compromise between FP and FN.  
 560 Moreover, it attains the lowest number of MCS. Therefore the proposed approach is well designed  
 561 to detect zero-day ransomware.

562

563 **Table 7:** Comparison of the proposed DCAE-ZSL-HVE with Deep Learning models in terms of  
 564 Recall, and, Accuracy on zero-day attacks based test data

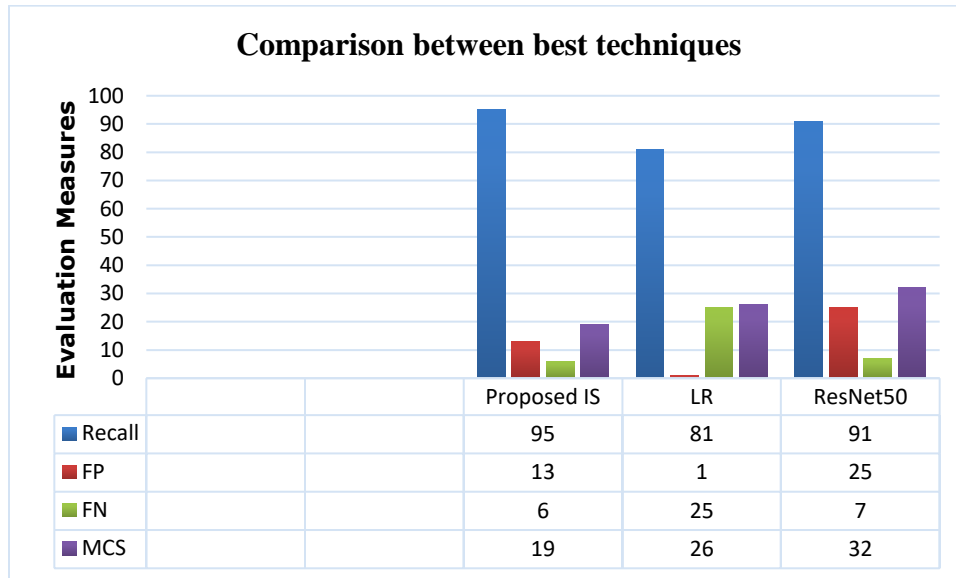
565

Techniques	TN	FP	FN	TP	Recall	Accuracy (%)
GoogleNet[46]	127	33	6	101	0.81	85.3
Inception V3[44]	127	23	6	111	0.84	89.1
ResNet50[43]	126	25	7	109	0.91	88.0
Proposed DCAE-ZSL-HVE (LM)	120	13	6	128	0.95	92.8

566

567

568  
569



570  
571  
572

**Fig. 11– Comparison of the proposed method with best performing state of the art method**

573 To evaluate the performance of the proposed methods against other detection systems, we  
574 investigate different approaches from the literature for ransomware detection. Table 8 shows the  
575 comparison between DCAE-ZSL and EldeRan [26], evaluated using ZSL-test-data given in  
576 section 3.3. Both techniques are compared in terms of the detection rate on 100 features. However,  
577 the proposed method is using less training data as compared to the existing technique. The  
578 proposed DCAE-ZSL technique does not include the remaining test classes in training while  
579 evaluating one family from ZSL-test data. The results show that the proposed technique is better  
580 in detecting Pgpocoder, Trojan-Ransom and Reveton than the existing method.

581  
582

**Table 8:** Detection rate comparison with existing EldeRan[26] on test data

Family	No of ransomware samples	Detection Rate (100 features)	
		EldeRan	DCAE-ZSL-HVE
<b>Pgpocoder</b>	4	0.75	<b>1</b>
<b>Reveton</b>	90	0.88	<b>0.96</b>
<b>TeslaCrypt</b>	6	1.00	0.92
<b>Trojan-Ransom</b>	34	0.94	<b>0.95</b>

583 Table 9 shows the comparison with the existing techniques in term of recall measure. We compared  
584 the recall performance of proposed DCAE-ZSL\_HVE (0.95) with some of the existing methods  
585 namely , C4.5(0.52)[13], KNN(0.70)[13], DNA Sequencing Engine(0.82) [47]and Anomaly(0.89)  
586 [27]. DNA Sequencing Engine is using the same dataset, but train test split may not be the same.

587 However, its testing is carried out on the known type of ransomware. Similarly, methods such as  
 588 C4.5[13] and KNN[13] are tested on known types of ransomware. Results show that the proposed  
 589 method is outperforming even on unknown test data compared to the existing methods that show  
 590 performance on known attacks. These results show that intra-family variations of ransomware can  
 591 be suppressed effectively using the proposed DCAE-ZSL technique. It also shows that by  
 592 efficiently integrating the decision of base learners trained on two different representation can  
 593 perform well than the single base learner.

594

595 **Table 9:** Recall comparison with existing techniques on test data

Techniques	Type	Recall	Test data
<b>Proposed DCAE-ZSL-HVE (Local Majority)</b>	Heterogenous Ensemble	<b>0.95</b>	Unknown attacks
<b>Saleh et al. [13]</b>	KNN	0.70	Known attacks
<b>Saleh et al. [13]</b>	C4.5	0.52	Known attacks
<b>Al-rimy et al.[27]</b>	Anomaly	0.89	Unknown attacks
<b>Sgandurra et al. [26]</b>	EldeRAN	<b>0.93</b>	UnKnown attacks

596

597 Table 10. shows the qualitative comparison using some significant characteristics between the  
 598 proposed method and the most recent ransomware detection methodologies. The quantitative  
 599 comparison is difficult to perform due to lack of availability of any benchmark dataset, different  
 600 test environments, and due to the difference in objectives of developing the detection systems. The  
 601 proposed method uses full set dynamic features to capture all malicious API sequences, file access,  
 602 registry access, drop files, file extension search, and threatened text instead of focusing on single  
 603 event. Further, we have extracted the core semantic embeddings that utilizes the subset of these  
 604 events. Further we objectively trained and tested the proposed ensemble model to detect zero-day  
 605 ransomware.

606 **Table 10:** Comparison of the proposed method with recent Ransomware detection techniques

Techniques	year	Ensemble	Features Type	Zero-day Attack Detection	Cost Sensitive Approach	Classification Type
<b>Proposed</b>	2021	yes	Entire dynamic features	Yes	yes	Binary
<b>B.Zhang et al.[48]</b>	2020	No	N-gram opcodes	No	No	Family
<b>Al-rimy et al.[49]</b>	2019	yes	Dynamic API Calls	No	No	Binary

<b>R.Vinayakumar et al.[50]</b>	2019	No	User centric data	No	No	Family
<b>S.Maniath et al.[51]</b>	2017	No	Dynamic API Calls	No	No	Binary
<b>S.Homayoun[52]</b>	2019	yes	Sequences of events	Yes	No	Family

## 6. Threats to validity

While developing the framework for zero-day ransomware detection, the proposed framework has some threats to validity which are defined below:

**Construct validity:** In the present research, the developed framework for ransomware detection only discriminate between the goodware and ransomware of local machine, but does not include the verification of network traffic with server. Moreover, family classification of ransomware is not performed.

**Internal validity:** Another threat is the homogeneity of the data used. However, in this current work, the data has been collected from diverse sources and thus it is largely heterogeneous.

**External Validity:** In this work, we considered 11 different ransomware families, while for training we used seven families and the remaining four families are used to evaluate the model for zero-day ransomware detection. Additionally, the work can be extended to train the framework with more ransomware families, which then can become more proficient to detect the real world ransomware. Another external threat to validity is timely detection, however dynamic analysis is a time taking process.

## 7. Conclusion and future work

This paper presents a two-stage ransomware detection system using concepts of ZSL, DCAE, and Ensemble learning. The first stage is an AL (Attribute Learning) phase, whereby, a novel DCAE-ZSL technique is proposed to learn the uniform latent semantic embeddings both for known and unknown attacks. The learned latent space also penalizes the input for little changes and can focus on major similarities between known and unknown classes. Thus, it is able to learn only the core transformation that results in performance improvement against zero-day attack detection. Through different experiments, it is observed that DCAE-ZSL based transformed representation outperforms the traditional machine learning approaches against zero-day attacks. Then, base learners are trained on the original features and on features extracted through the proposed DCAE-ZSL method. Finally, the second stage of the proposed framework consists of the IS (Inference Stage) that implements four combination rules to obtain the final prediction results. The prediction results by K-NN, DNA Sequencing Engine and Anomaly based methods shows that our method

637 can provide good recall score as compared to these methods. Further, the proposed method shows  
638 better performance in detecting Pgpocoder, Trojan-Ransom, and Reveton families than the existing  
639 method [20] in terms of precision. Moreover, results show that the proposed DCAE-ZSL-HVE  
640 methods achieve a considerable compromise between FP and FN as compared to the conventional  
641 baseline models. Hence, we conclude that the invariant and reduced feature representation of the  
642 original features can efficiently detect new classes at the test time.  
643 In this study, the developed system for zero-day ransomware detection only detects either an exe  
644 file is a goodware or ransomware. In future, this work can be extended to identify the family of  
645 the ransomware. Similarly, in future CAE can also be used to extract the pre- encryption based  
646 discriminant features for an early detection.

---

## 647 **Acknowledgment**

648 The authors would like to thank Pattern Recognition lab at Pakistan Institute of Engineering and  
649 Applied Sciences for providing them healthy research environment and computational facilities.  
650

651

652 **References**

- 653 1. Al-rimy BAS, Maarof MA, Shaid SZM (2018) Ransomware threat success  
654 factors, taxonomy, and countermeasures: A survey and research directions.  
655 Comput. Secur.
- 656 2. Hwang J, Kim J, Lee S, Kim K (2020) Two-Stage Ransomware Detection  
657 Using Dynamic Analysis and Machine Learning Techniques. Wirel Pers  
658 Commun. <https://doi.org/10.1007/s11277-020-07166-9>
- 659 3. Bilge L, Dumitras T (2012) Before we knew it: An empirical study of zero-  
660 day attacks in the real world. In: Proceedings of the ACM Conference on  
661 Computer and Communications Security
- 662 4. Sood AK, Enbody RJ (2011) Malvertising - Exploiting web advertising.  
663 Comput Fraud Secur. [https://doi.org/10.1016/S1361-3723\(11\)70041-0](https://doi.org/10.1016/S1361-3723(11)70041-0)
- 664 5. McDowell K (2006) Now that we are all so well-educated about spyware, can  
665 we put the bad guys out of business? In: Proceedings of the 34th Annual ACM  
666 SIGUCCS Fall 2006 Conference, SIGUCCS '06
- 667 6. Kharraz A, Robertson W, Balzarotti D, et al (2015) Cutting the gordian knot:  
668 A look under the hood of ransomware attacks. In: Lecture Notes in Computer  
669 Science (including subseries Lecture Notes in Artificial Intelligence and  
670 Lecture Notes in Bioinformatics)
- 671 7. Masdari M, Khezri H (2020) A survey and taxonomy of the fuzzy signature-  
672 based Intrusion Detection Systems. Appl. Soft Comput. J.
- 673 8. Medhat M, Gaber S, Abdelbaki N (2018) A new static-based framework for  
674 ransomware detection. In: Proceedings - IEEE 16th International Conference  
675 on Dependable, Autonomic and Secure Computing, IEEE 16th International  
676 Conference on Pervasive Intelligence and Computing, IEEE 4th International  
677 Conference on Big Data Intelligence and Computing and IEEE 3
- 678 9. Jung S, Won Y (2018) Ransomware detection method based on context-aware  
679 entropy analysis. Soft Comput. <https://doi.org/10.1007/s00500-018-3257-z>
- 680 10. Zakaria WZA, Abdollah MF, Mohd O, Ariffin AFM (2017) The rise of  
681 ransomware. In: ACM International Conference Proceeding Series
- 682 11. Lin SW, Ying KC, Lee CY, Lee ZJ (2012) An intelligent algorithm with  
683 feature selection and decision rules applied to anomaly intrusion detection.  
684 Appl Soft Comput J. <https://doi.org/10.1016/j.asoc.2012.05.004>
- 685 12. Tajoddin A, Abadi M (2019) RAMD: registry-based anomaly malware

686 detection using one-class ensemble classifiers. Appl Intell.  
687 <https://doi.org/10.1007/s10489-018-01405-0>

688 13. Saleh MA (2020) Evaluation of Supervised Machine Learning Classifiers for  
689 Detecting Ransomware based on Naïve Bayes, SVM, KNN, C 4.5,  
690 andRandom Forest Algorithms

691 14. Ashit D (2016) Detection of Malware and Malicious Executables Using E-  
692 Birch Algorithm. Int J Adv Comput Sci Appl.  
693 <https://doi.org/10.14569/ijacsa.2016.070118>

694 15. Tian Q, Han D, Li KC, et al (2020) An intrusion detection approach based on  
695 improved deep belief network. Appl Intell. <https://doi.org/10.1007/s10489-020-01694-4>

696

697 16. Idhammad M, Afdel K, Belouch M (2018) Semi-supervised machine learning  
698 approach for DDoS detection. Appl Intell. <https://doi.org/10.1007/s10489-018-1141-2>

699

700 17. Andronio N, Zanero S, Maggi F (2015) HELDROID: Dissecting and detecting  
701 mobile ransomware. In: Lecture Notes in Computer Science (including  
702 subseries Lecture Notes in Artificial Intelligence and Lecture Notes in  
703 Bioinformatics)

704 18. Mercaldo F, Nardone V, Santone A, Visaggio CA (2016) Ransomware steals  
705 your phone. Formal methods rescue it. Lect Notes Comput Sci (including  
706 Subser Lect Notes Artif Intell Lect Notes Bioinformatics).  
707 [https://doi.org/10.1007/978-3-319-39570-8\\_14](https://doi.org/10.1007/978-3-319-39570-8_14)

708 19. Das S, Liu Y, Zhang W, Chandramohan M (2016) Semantics-based online  
709 malware detection: Towards efficient real-time protection against malware.  
710 IEEE Trans Inf Forensics Secur. <https://doi.org/10.1109/TIFS.2015.2491300>

711 20. Alsoghyer S, Almomani I (2019) Ransomware detection system for android  
712 applications. Electron. <https://doi.org/10.3390/electronics8080868>

713 21. Kharraz A, Arshad S, Mulliner C, et al (2016) Unveil: A large-scale, automated  
714 approach to detecting ransomware. In: Proceedings of the 25th USENIX  
715 Security Symposium

716 22. Song S, Kim B, Lee S (2016) The Effective Ransomware Prevention  
717 Technique Using Process Monitoring on Android Platform. Mob Inf Syst.  
718 <https://doi.org/10.1155/2016/2946735>

719 23. Aydin MA, Zaim AH, Ceylan KG (2009) A hybrid intrusion detection system  
720 design for computer network security. Comput Electr Eng.

721 <https://doi.org/10.1016/j.compeleceng.2008.12.005>

722 24. Ferrante A, Malek M, Martinelli F, et al (2018) Extinguishing ransomware - a  
723 hybrid approach to android ransomware detection. In: Lecture Notes in  
724 Computer Science (including subseries Lecture Notes in Artificial Intelligence  
725 and Lecture Notes in Bioinformatics)

726 25. Netto DF, Shony KM, Lalson ER (2018) An Integrated Approach for  
727 Detecting Ransomware Using Static and Dynamic Analysis. In: 2018  
728 International CET Conference on Control, Communication, and Computing,  
729 IC4 2018

730 26. Sgandurra D, Muñoz-González L, Mohsen R, Lupu EC (2016) Automated  
731 Dynamic Analysis of Ransomware: Benefits, Limitations and use for  
732 Detection

733 27. Al-rimy BAS, Maarof MA, Prasetyo YA, et al (2018) Zero-Day Aware  
734 Decision Fusion-Based Model for Crypto-Ransomware Early Detection. *Int J*  
735 *Integr Eng* 10:. <https://doi.org/10.30880/ijie.2018.10.06.011>

736 28. Al-rimy BAS, Maarof MA, Shaid SZM (2018) A 0-Day Aware Crypto-  
737 Ransomware Early Behavioral Detection Framework

738 29. Sun X, Gu J, Sun H (2020) Research progress of zero-shot learning. *Appl Intell*

739 30. Rivero J, Ribeiro B, Chen N, Leite FS (2017) A grassmannian approach to  
740 zero-shot learning for network intrusion detection. In: Lecture Notes in  
741 Computer Science (including subseries Lecture Notes in Artificial Intelligence  
742 and Lecture Notes in Bioinformatics)

743 31. Zhang X, Gao L, Jiang Y, et al (2019) A zero-shot intrusion detection method  
744 based on regression model. *Proc - 2019 7th Int Conf Adv Cloud Big Data,*  
745 *CBD 2019* 186–191. <https://doi.org/10.1109/CBD.2019.00042>

746 32. Li Z, Qin Z, Shen P, Jiang L (2019) Zero-shot learning for intrusion detection  
747 via attribute representation. In: Lecture Notes in Computer Science (including  
748 subseries Lecture Notes in Artificial Intelligence and Lecture Notes in  
749 Bioinformatics)

750 33. Lv N, Chen C, Qiu T, Sangaiah AK (2018) Deep Learning and Superpixel  
751 Feature Extraction Based on Contractive Autoencoder for Change Detection  
752 in SAR Images. *IEEE Trans Ind Informatics*.  
753 <https://doi.org/10.1109/TII.2018.2873492>

754 34. Van Der Maaten L, Hinton G (2008) Visualizing data using t-SNE. *J Mach*  
755 *Learn Res*

- 756 35. Baldi P (2012) Autoencoders, Unsupervised Learning, and Deep  
757 Architectures. ICML Unsupervised Transf Learn.  
758 <https://doi.org/10.1561/22000000006>
- 759 36. Qureshi AS, Khan A, Shamim N, Durad MH (2019) Intrusion detection using  
760 deep sparse auto-encoder and self-taught learning. Neural Comput Appl.  
761 <https://doi.org/10.1007/s00521-019-04152-6>
- 762 37. Rifai S, Vincent P, Muller X, et al (2011) Contractive auto-encoders: Explicit  
763 invariance during feature extraction. In: Proceedings of the 28th International  
764 Conference on Machine Learning, ICML 2011
- 765 38. Romera-Paredes B, Torr PHS (2015) An embarrassingly simple approach to  
766 zero-shot learning. In: 32nd International Conference on Machine Learning,  
767 ICML 2015
- 768 39. Lampert CH, Nickisch H, Harmeling S (2014) Attribute-based classification  
769 for zero-shot visual object categorization. IEEE Trans Pattern Anal Mach  
770 Intell. <https://doi.org/10.1109/TPAMI.2013.140>
- 771 40. Akata Z, Perronnin F, Harchaoui Z, Schmid C (2013) Label-embedding for  
772 attribute-based classification. In: Proceedings of the IEEE Computer Society  
773 Conference on Computer Vision and Pattern Recognition
- 774 41. Lampert CH, Nickisch H, Harmeling S (2010) Learning to detect unseen object  
775 classes by between-class attribute transfer
- 776 42. Binbusayyis A, Vaiyapuri T (2021) Unsupervised deep learning approach for  
777 network intrusion detection combining convolutional autoencoder and one-  
778 class SVM. Appl Intell
- 779 43. He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image  
780 recognition. In: Proceedings of the IEEE Computer Society Conference on  
781 Computer Vision and Pattern Recognition
- 782 44. Szegedy C, Vanhoucke V, Ioffe S, et al (2016) Rethinking the Inception  
783 Architecture for Computer Vision. In: Proceedings of the IEEE Computer  
784 Society Conference on Computer Vision and Pattern Recognition
- 785 45. Vasan D, Alazab M, Wassan S, et al (2020) Image-Based malware  
786 classification using ensemble of CNN architectures (IMCEC). Comput Secur.  
787 <https://doi.org/10.1016/j.cose.2020.101748>
- 788 46. Zeng G, He Y, Yu Z, et al (2016) InceptionNet/GoogLeNet - Going Deeper  
789 with Convolutions. Cvpr
- 790 47. Khan F, Ncube C, Ramasamy LK, et al (2020) A Digital DNA Sequencing

791 Engine for Ransomware Detection Using Machine Learning. IEEE Access.  
792 <https://doi.org/10.1109/ACCESS.2020.3003785>

793 48. Zhang B, Xiao W, Xiao X, et al (2020) Ransomware classification using patch-  
794 based CNN and self-attention network on embedded N-grams of opcodes.  
795 *Futur Gener Comput Syst.* <https://doi.org/10.1016/j.future.2019.09.025>

796 49. Al-rimy BAS, Maarof MA, Shaid SZM (2019) Crypto-ransomware early  
797 detection model using novel incremental bagging with enhanced semi-random  
798 subspace selection. *Futur Gener Comput Syst.*  
799 <https://doi.org/10.1016/j.future.2019.06.005>

800 50. Vinayakumar V, Alazab M, Jolfaei A, et al (2019) Ransomware triage using  
801 deep learning: Twitter as a case study. In: *Proceedings - 2019 Cybersecurity  
802 and Cyberforensics Conference, CCC 2019*

803 51. Maniath S, Ashok A, Poornachandran P, et al (2018) Deep learning LSTM  
804 based ransomware detection. In: *2017 Recent Developments in Control,  
805 Automation and Power Engineering, RDCAPE 2017*

806 52. Homayoun S, Dehghantanha A, Ahmadzadeh M, et al (2019) DRTHIS: Deep  
807 ransomware threat hunting and intelligence system at the fog layer. *Futur  
808 Gener Comput Syst.* <https://doi.org/10.1016/j.future.2018.07.045>

809

#### 810 **Conflicts of interest**

811 Authors declare no conflict of interest.

#### 812 **Availability of data and material**

813 All the datasets used in this work are publicly available, whereas datasets that are generated during  
814 simulations are available from the corresponding author on reasonable request.

815