# Malware Classification Using Deep Learning Adversarial Attack Images Detection

## Sri Lanka Institute of Information Technology

### Individual Assignment

IE4032 – Information Warfare

Submitted by:

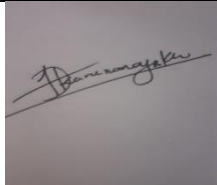| Student Registration Number | Student Name |
|---|---|
| IT20665784 | H.D.Karunanayaka |

Date of submission

2023/11/03

# DECLARATION

I declare that this is my own work, and this proposal does not incorporate without acknowledgement any material previously submitted for a degree or diploma in any other university or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

**Signature:**

| IT20665784 | H.D.Karunanayaka | | |
|---|---|---|---|
| | | | |

The above candidates are carrying out research for the undergraduate Dissertation under my supervision.

## ACKNOWLEDGEMENTS

## ABSTRACT

Adversarial attacks exploit the vulnerability of deep neural network (DNN) models by introducing imperceptibly small yet carefully crafted attack data, causing significant problems. In response to these attacks, various defense strategies have been proposed: (1) providing adversarial training based on specific attacks, (2) reducing input data noise, (3) performing input data preprocessing, and (4) incorporating noise into multiple layers of the model. In this context, we present a simple yet effective approach called the Noise-Fusion Method (NFM) to combat adversarial attacks on DNN image classification models. The NFM not only applies noise to the model input during runtime but also incorporates noise into the training data during the training process, without requiring knowledge about specific attacks or models. To evaluate the defensive capabilities of the NFM, we employ the Sparse L1 Descent (SLD) method as well as two popular adversarial attacks, namely the Fast Gradient Signed Method (FGSM) and the Projected Gradient Descent (PGD). We conduct experiments using the MNIST and CIFAR-10 datasets with various deep neural network models.

To demonstrate the effectiveness of the NFM against different types of noise, we introduce a range of amplitude sounds with diverse statistical distributions. Our results indicate that the inclusion of noise in both the input and training images not only enhances the resilience of the associated models but also safeguards against all three adversarial attacks.

Keywords—adversarial attack; defense; noise; deep neural network

**Table Of Content**

# 1.0 INTRODUCTION OF PRODUCT

The "Malware Classification Using Deep Learning and Adversarial Attack Images Detection" Power BI dashboard is a powerful and innovative tool designed to address the growing cybersecurity challenges in today's digital landscape. This dashboard serves as a visual representation of our cutting-edge research and methodology for malware classification and detection.

In a world where cyber threats and malware attacks continue to evolve at an alarming rate, it is essential to have robust defense mechanisms in place. Traditional signature-based approaches are no longer sufficient to protect against advanced and adaptive malware. This is where deep learning and adversarial attack detection play a pivotal role.

Our dashboard provides a comprehensive and intuitive interface for security professionals, analysts, and decision-makers to gain valuable insights into the world of malware and cyber threats. It offers a visually appealing and user-friendly platform that empowers users to:

1. **Monitor Malware Trends**: Stay updated on the latest malware trends and types through data visualization and analysis.

2. **Classify Malware**: Understand the classification of malware into different categories using our deep learning models.

3. **Detect Adversarial Attacks**: Identify and respond to adversarial attacks targeting machine learning models for malware detection.

4. **Analyze Model Performance**: Evaluate the effectiveness of our deep learning models in real-time.

5. **Enhance Cybersecurity**: Make informed decisions to enhance cybersecurity strategies and protect digital assets.

This dashboard is a result of extensive research and development, and it offers a comprehensive solution to address the complex challenges posed by malware and

adversarial attacks. It provides valuable data-driven insights that help organizations stay ahead of cyber threats and secure their systems and data effectively.

In the following sections of this report, we will delve deeper into the scope, methodology, and key findings of our research to shed light on the significance and functionality of this Power BI dashboard.
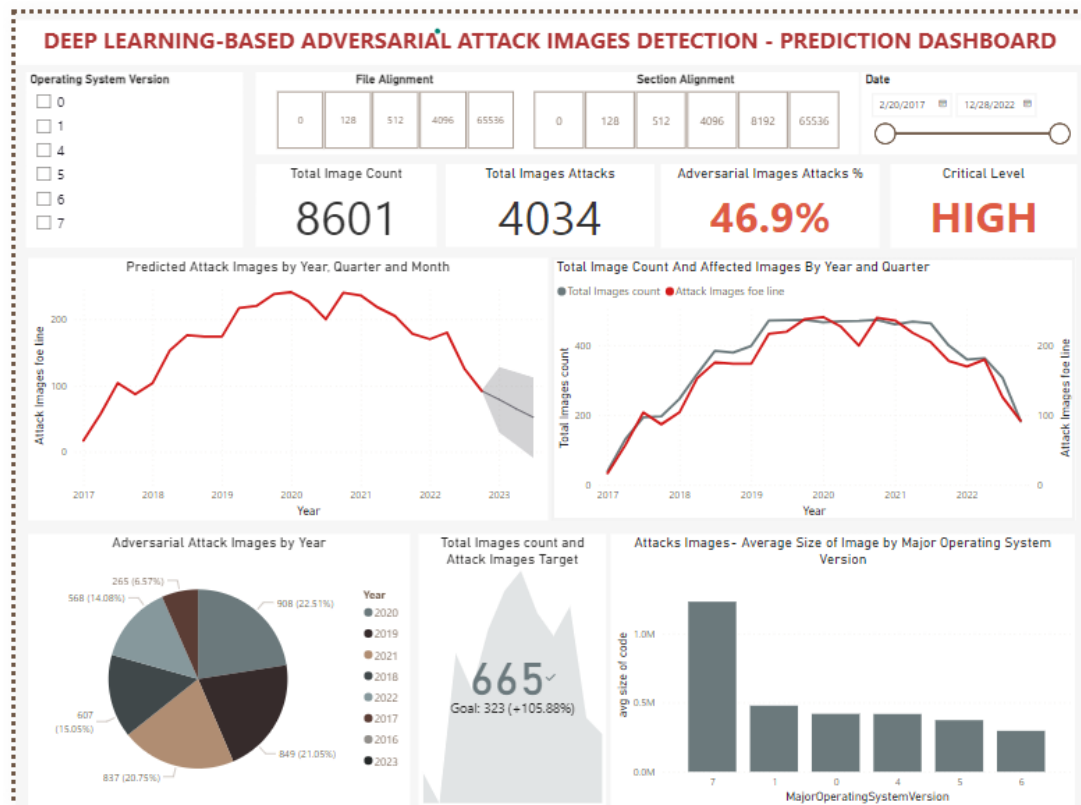


Fig 1: Deep learning based adversarial attack images detection – Power BI Dashboard.

First developed by Krizgevsky et al. in 2012, the Convolutional Neural Network (CNN) was the top model in a difficult image classification competition using the ImageNet dataset. Since then, one of the most well-known machine learning technologies is deep learning, which has grown quickly. Deep learning-based artificial intelligence techniques are widely used in vital domains including speech recognition, natural language processing, and picture classification. However, recent research has revealed that adversarial examples can compromise and undermine the

effectiveness of deep neural network models. Even a slight perturbation added to a regular input image can cause the model's output to deviate from the correct prediction, progressively decreasing its classification accuracy. According to Szegedy et al., even imperceptible alterations to an image can lead to significant changes in the prediction of a deep neural network, often invisible to the human visual system. For instance, an image initially classified as a "panda" can be misclassified as a "gibbon" after a minor perturbation. Moreover, attacked models exhibit high confidence in these deceivingly altered images, and the same perturbation can fool multiple similar models. These findings demonstrate the considerable harm that adversarial examples pose to machine learning, particularly in deep learning scenarios.
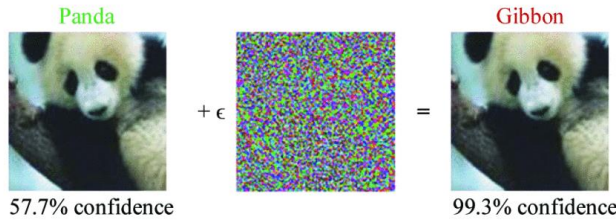


Fig 2. setting an antagonistic example. The neural network's input image is shown on the left, a rapidly generated little perturbation is shown in the middle, and an adversarial example with an additional disturbance is shown on the right.

Adversarial examples are not limited to the digital realm; they can also be used to exploit deep learning models in the physical world. For instance, fake road signs could mislead an autonomous vehicle into making mistakes. Real-time applications of deep neural networks, such as face recognition and 3D object identification, also face similar threats. We concentrate on traditional white-box and black-box attack techniques in this study since they have been widely explored and optimized. White-box approaches are frequently used to assess the model's defensive performance and robustness. Additionally, we provide a thorough examination of defensive mechanisms used against hostile instances, classifying them according to their source of defense efficacy. We test alternative defensive strategies using the same dataset and assault settings because there is no defined contrast mechanism. While adversarial examples exist across a range of deep learning tasks, research on adversarial examples in the image domain is more extensive. Therefore, this work primarily focuses on adversarial attacks and defenses in the context of image

categorization tasks. The subsequent sections of this paper are organized as follows: Section II outlines the approaches for generating adversarial examples. Section III presents and elaborates on defense approaches from four perspectives: gradient masking, adversarial training, adversarial example detection, and input modification. Section IV provides a comprehensive review of the performance and characteristics of defense strategies against adversarial examples. Furthermore, We explore the research issues in the realm of adversarial example defenses and provide insights into future research opportunities.

attack strategies in this study since there has been a lot of important work done to optimize and improve these approaches. Furthermore, white-box approaches are commonly used to assess the model's defense performance or robustness. We also perform a thorough investigation of the defensive mechanisms employed in adversarial instances and categorize them depending on the source of the defense impact. Because there is no established contrast mechanism in the current work, we attempt to evaluate various defense methods against examples using the identical dataset and attack parameters in this paper. Although adversarial examples can be found in deep learning models for a variety of tasks, research on adversarial in this study, we aim to assess alternative defensive strategies against instances using the same dataset and assault settings because there is no proven contrast mechanism in the current work. Although hostile examples may be found in deep learning models for a wide range of tasks, research on adversarial examples in images is more extensive. As a result, this work focuses primarily on adversarial assaults and defenses of picture categorization tasks. The rest of this paper is arranged as follows. Section II describes the approaches for generating hostile instances. Section III provides and elaborates on adversarial example defense approaches from four perspectives: gradient masking, adversarial training, adversarial examples detection, and input modification. Section IV reviews the performance and features of hostile instances of defense strategies in depth. Furthermore, we discuss the research challenges in the field of adversarial examples of defense techniques, as well as the future research prospects.

# I. ADVERSARIAL ATTACK

Adversarial instances take advantage of flaws in deep neural network models by purposely changing input examples to fool the networks. Minor alterations to the picture are introduced during image classification to render the categorization label erroneous. Adversarial assaults are classed as either target or non-target based on their intended consequence. The goal of targeted assaults is to modify the model's categorization to a specific label that is usually different from the initial label. Non-target assaults simply require a categorization that differs from the original label. This article largely focuses on target assaults and their related responses.

Based on the attacker's information, adversarial assaults are divided into three types: black-box attacks, grey-box attacks, and white-box attacks. The attacker in a black-box scenario is unaware of the model's parameters, training dataset distribution, or protection methods. The attacker in the grey-box scenario understands the model's parameters and the distribution of the training dataset but is unaware of the precise protection mechanisms. In the white-box scenario, the attacker has access to all of the model's parameters, as well as the distribution of training datasets and protection mechanisms. This work explores both black-box and white-box attack strategies.

Adversarial attack approaches can be further classified based on the creation process of adversarial examples, including gradient-based attacks, optimization-based attacks, score-based attacks, and decision-based attacks. These categories are determined by the source of the attack's impact.

In this research paper, the focus is on adversarial attacks, which are techniques used to exploit the weaknesses of deep neural network models by intentionally manipulating input examples. Adversarial attacks aim to create adversarial examples, which are modified versions of the original input that can deceive the neural network into producing incorrect or misleading outputs.

The paper discusses various types of adversarial attacks based on different factors. One categorization is based on the desired outcome of the attack: target attacks and non-target attacks. Target attacks aim to misclassify the input as a specific target class, while non-target attacks seek to change the classification without a specific target in mind.

Another classification considered in the paper is based on the attacker's knowledge and access to information. Adversarial attacks can be categorized as black-box attacks, grey-box attacks, or white-box attacks. In black-box attacks, the attacker has no knowledge of the

model's parameters, training dataset distribution, or defense mechanisms. Grey-box attacks involve partial knowledge, such as knowing the model's parameters and training dataset distribution but not the specific defense measures. White-box attacks grant the attacker full access to the model's parameters, training dataset distribution, and defense techniques.

The paper also discusses different approaches for generating adversarial examples. These approaches include gradient-based attacks, optimization-based attacks, score-based attacks, and decision-based attacks. Each approach has its own methodology for manipulating the input examples to create adversarial perturbations that can deceive the neural network.

Overall, this section of the research paper provides an overview of adversarial attacks, their categorization based on various factors, and the approaches used to generate adversarial examples. Understanding adversarial attacks is crucial for developing robust defense mechanisms to protect deep neural network models against such threats.

*A.* Gradient-based Attacks

Gradient-based attack techniques involve backpropagating through the neural network's cost function to compute gradients with respect to the input image. These gradients are then utilized to modify the input in order to create an adversarial example. One of the commonly used methods in this category is the Fast Gradient Sign Method (FGSM), introduced by Goodfellow et al. FGSM manipulates the image by maximizing the loss of the classifier on the image. Mathematically, an adversarial example is generated using the following equation, where $\varepsilon$ denotes the perturbation magnitude, sign() extracts the sign of the gradient, and J() represents the cost function used for training the model.

$$x' = x - \epsilon \cdot sign\left(\nabla_x J_\theta(F(x))\right)$$

FGSM is efficient as it requires only a single computation step to generate adversarial samples. Additionally, Goodfellow et al. proposed the Least-Likely-Class (LLC) attack as an extension of FGSM, where the target label with the lowest class probability replaces the original class variable. LLC is a targeted variant of the FGSM attack.

The Basic Iterative Method (BIM) is another extension of FGSM that perturbs the image iteratively over multiple small steps, adjusting the gradient direction after each step to improve effectiveness. Momentum concepts were further integrated into BIM, leading to the creation of the Momentum Iterative Fast Gradient Sign Method (MI-FGSM). MI-FGSM improved the stability of gradient direction updates and addressed the issue of BIM getting trapped in local minima when generating adversarial examples.

$$x'_i = x'_{i-1} - clip_\epsilon\left(\alpha \cdot sign\left(\nabla_x J_\theta(F(x'_{i-1}))\right)\right)$$

The maximum distortion is controlled by xi1 ′, where xi1 ′ is the updated input at step i-1. PGD was the most powerful first-order attack technique at the time, and many advanced attack methods are enhanced variants of PGD. This assault approach is rapid and mobile.

The Projected Gradient Descent (PGD) approach, which uses randomly initialized projected gradient descent, was developed as an extension to BIM. PGD perturbs the input picture several times with a fixed step size based on the gradient's sign direction. The maximum distortion is controlled by the variable xi1' in the updated sample at step i-1. PGD was regarded as one of the most potent first-order assault methods, and many modern attack methods are upgraded variations of PGD.

In summary, gradient-based attack methods are rapid and mobile, involving manipulation of the input image using gradients derived from the cost function. FGSM, LLC, BIM, MI-FGSM, and PGD are notable examples in this category, with PGD being a particularly powerful and widely used technique.

### B. Optimization-based Attacks

Optimization-based assaults seek adversarial instances by resolving the problem as an optimization work with the goal of locating hostile examples with high confidence and low disturbance. Szegedy et al. proposed the Limited Memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) attack, which proved that introducing undetectable disruptions to an image may cause neural networks to misclassify it. Using the optimization equation:

$$min \quad c\|x - x'\| + L(\theta, x', t), \quad s.t. \quad x' \in [0,1]^m$$

In the L-BFGS attack, the goal is to minimize the classifier loss L() while introducing perturbations to the input image. However, this approach has the drawback of being computationally expensive and impractical due to the need for a large-scale linear search to find the optimal solution.

To quantify the magnitude of the perturbation, various norm measurements are employed, leading to different types of attacks such as L0 attack, L2 attack, and L attack. These measurements help assess the level of disturbance introduced to the image. Notably, optimization-based attacks can often achieve a high success rate with minimal noise, making them particularly effective in certain scenarios.

Overall, optimization-based attacks leverage well-designed optimization techniques to find adversarial examples by minimizing the classifier loss while introducing minimal perturbations. The L-BFGS attack is an example of this approach, although its computational requirements limit its practicality. Different norm measurements are used to quantify the perturbation magnitude, allowing for fine-tuning of the attack based on the desired level of disturbance.

*C.* PGD (Projected Gradient Descent) Attack

The Projected Gradient Descent (PGD) attack is a robust adversarial technique widely employed to assess the resilience of machine learning models. In PGD, the attacker calculates the gradient of a loss function with respect to the input and then adjusts the input in the direction of the gradient that maximizes the loss. To ensure that the modified input remains within the valid input space for the model, the attacker projects it onto a permissible set, typically defined by constraints such as a range of pixel values.

By iteratively applying this method, the attacker generates adversarial examples that appear visually indistinguishable from the original input but cause the model to confidently misclassify them. PGD is a popular choice for both research and practical applications since it is both effective and adaptable. It may be used to attack several types of machine learning models, including deep neural networks, and can be tailored to different threat models, such as black-box assaults, in which the attacker has limited knowledge of the target model.

PGD's ability to generate adversarial examples that retain perceptual similarity to the original input while fooling the model makes it a valuable tool for evaluating model robustness and investigating potential vulnerabilities. Its flexibility and broad applicability contribute to its significance in adversarial research and highlight its importance in understanding and improving the security of machine learning systems.

*D.* DEFENSE AGAINST ADVERSARIAL ATTACKS

We investigate the work of representative adversarial defenses such as gradient masking, adversarial training, detection, and input alterations in this paper. In addition, we introduce several unique protection techniques.

a. Gradient Masking

Gradient masking is a defensive approach that conceals gradients in order to reduce a model's susceptibility to hostile samples. Paper not et al. adopted a strategy known as defensive distillation, which reduces the amplitude of network gradients to make the model less vulnerable to perturbations induced by adversarial cases. Ross and Doshi-Velez discovered, however, that the hardened model with defensive distillation did not outperform the traditional unprotected model. As a result, they developed the input gradient regularization technique, which directly improves the model to have smoother input gradients, hence increasing the model's robustness. This method trains a differentiable model to punish both input and output changes, guaranteeing that even minor changes do not have a substantial impact on the model's output.

When combined with adversarial training, this approach demonstrates a positive impact on the model's resilience against attacks such as FGSM (Fast Gradient Sign Method) and JSMA (Jacobian-based Saliency Map Attack). However, one drawback of this method is its high training complexity. To enhance the performance of gradient regularization in classification tasks, Yeats et al. proposed the use of

complex-valued neural networks (CVNNs). The gradient masking defense technique primarily focuses on disguising the initial model gradients, making it highly effective against adversarial examples generated using gradient-based methods. However, its defense effectiveness decreases when faced with adversarial examples generated by other techniques.

In summary, gradient masking is a defense mechanism that aims to reduce a model's vulnerability to adversarial examples by concealing the gradients. Defensive distillation and input gradient regularization are two approaches used in this context, with the latter showing promising results when combined with adversarial training.

However, the high training complexity and limited effectiveness against certain types of adversarial examples are important considerations when employing gradient masking as a defense strategy.

b. Adversarial Training

| Taxonomy | Paper | Model Architecture | Attack | $\epsilon$ | Dataset | Accuracy |
|---|---|---|---|---|---|---|
| Adversarial Regularization | [40] | ResNet-152 | $PGD_{50}$ | 4/255 | ImageNet | 47.00% |
| | [41] | Wide ResNet | $CW_{10}$ | 0.031/1 | CIFAR-10 | 84.03% |
| | [42] | ResNet-18 | $PGD_{20}$ | 8/255 | CIFAR-10 | 55.45% |
| | [43] | InceptionV3 | $PGD_{20}$ | 16/255 | ImageNet | 27.90% |
| Curriculum | [44] | Wide ResNet | $PGD_{20}$ | 16/255 | CIFAR-10 | 49.86% |
| | [45] | DenseNet-161 | $PGD_7$ | 8/255 | CIFAR-10 | 69.27% |
| | [46] | 8-Layer ConvNet | $PGD_{20}$ | 8/255 | CIFAR-10 | 42.40% |
| Ensemble | [47] | Wide ResNet | $PGD_{10}$ | 0.005 | CIFAR-100 | 32.10% |
| | [48] | ResNet-20 | $PGD_{50}$ | 0.09/1 | CIFAR-10 | 46.30% |
| | [49] | ResNet-20 | PGD20 | 0.01/1 | CIFAR-10 | 52.4% |
| Adaptive $\epsilon$ | [50] | ResNet-152 | $PGD_{1000}$ | 8/255 | ImageNet | 59.28% |
| | [51] | Wide ResNet | $PGD_{100}$ | 8/255 | CIFAR-10 | 47.18% |
| | [52] | Wide ResNet | $PGD_{20}$ | 8/255 | CIFAR-10 | 73.38% |
| Semi-Unsupervised | [53] | Wide ResNet | FGSM | 8/255 | CIFAR-10 | 62.18% |
| | [54] | Wide ResNet | PGD10 | 8/255 | CIFAR-10 | 63.10% |
| | [55] | Wide ResNet | PGD20 | 0.3/1 | ImageNet | 50.40% |
| Efficient | [56] | Wide ResNet | PGD100 | 8/255 | CIFAR-10 | 46.19% |
| | [57] | ResNet-50 | PGD40 | 2/255 | ImageNet | 43.43% |
| | [58] | PreActResNet-18 | FGSM | 8/255 | CIFAR-10 | 50.50% |
| Benchmark | [18] | ResNet-50 | PGD20 | 8/255 | CIFAR-10 | 45.80% |

TABLE 1 An overview of the experimental outcomes for several adversarial training approaches.

Goodfellow et al. pioneered the notion of adversarial training. The primary concept behind this method is to include adversarial instances in the training set and train the model using both actual and adversarial examples. During training, the model can learn how to handle hostile instances. Researchers have devised several approaches to increase the model's resilience using adversarial instances since its introduction.

Goodfellow et al. created adversarial examples for training data using the FGSM technique. Kurakin et al. improved on this strategy by including batch standardization, resulting in a mixed training approach with superior experimental outcomes. However, subsequent research showed that FGSM adversarial training was not always effective in improving the model's robustness and could reduce prediction accuracy in the face of more powerful iterative attacks. As a result, Madry et al. proposed using stronger PGD attacks for adversarial training and formulated

the process as an optimization problem of minimization and maximization. Building on this approach, Kan proposed logit pairing to further enhance robustness.

Adversarial examples are transferable, which means that creating adversarial examples for one model can lead to inaccurate predictions for another. . Kurakin et al. solved this issue by introducing an integrated adversarial training strategy that created adversarial instances from several pre-training models and utilized them to enlarge the original data set.

Despite being considered the best defense method currently available, the participation of adversarial examples during training can significantly prolong the training time

The disadvantage of adversarial training is that it reduces the model's accuracy in predicting normal cases. Zhang et al. [41] presented the TRADES method to manage the trade-off between model robustness and accuracy. It smooths the decision boundary and reduces the prediction gap between actual and adversarial samples by adding a regularization loss to the loss function. Wang et al. [59] enhanced adversarial training even further by leveraging feature dispersion to produce adversarial instances with the greatest feature matching gap between genuine and adversarial examples. To improve the network's resilience, Zoran et al. [63] presented an attention model based on human perception. The model employs an attention mechanism to mimic the human visual system, and the picture is examined through a succession of glances. Zhang described the susceptibility of adversarial instances to deep neural networks in terms of neuron sensitivity and presented the SNS defensive approach, which stabilizes the differences between sensitive neurons in genuine and adversarial examples to increase adversarial resilience.

Bai et al. have offered a summary of recent achievements in adversarial training, categorizing them by distinct understandings of adversarial training. TABLE 1

highlights various adversarial training approaches, such as adversarial regularization, curriculum, ensemble training, adaptive, semi-unsupervised learning, efficient training, and others, as paths for further enhancement of adversarial training.

c.  Detection

The transformation-based strategy is based on the instability of adversarial instances. By altering the example before feeding it to the model for classification, the possible adversarial perturbation is removed. Dziugaite et al. [72] discovered that the majority of training photos are in JPG format and sought to reduce the influence of adversarial perturbations on accuracy by using JPG image compression. However, experiments showed that this method could only eliminate small perturbations and did not effectively address larger ones, resulting in reduced accuracy. Similarly, the proposed PCA compression method [73] suffered from the same limitations. Instead of using standard image compression techniques, there is literature suggesting the use of a learnable compression method employing neural network models. Theagarajan et al. [74] developed a defense module that learned to project inputs onto the non-adversarial data region of the target model. In a related work, Sun et al. [75] utilized convolutional sparse coding to transform input images. By employing a "sparse transformation layer" to project inputs into a quasi-natural space, the authors claimed reduced sensitivity to adversarial attacks. Xie et al. [76] discovered that random rescaling of adversarial examples weakened the intensity of attacks, while random image padding reduced the chances of fooling the network. Building on image compression, Guo et al. [77] proposed a method based on total variation minimization to eliminate perturbations. Song et al. [78] introduced the Pixeldefend method, which aimed to restore a maliciously perturbed image to a distribution resembling the training data. Brama et al. [79] proposed a defense method that combined feature visualization with input modification, making it applicable to various pretrained networks. They found that the heat maps generated by the LRP method exhibited minimal changes between real and adversarial examples. As a result, they blurred the input image and reconstructed it by adding relevant pixels obtained from LRP interpretation. The reconstructed image was then used as the model's input for re-prediction. It is important to note that general input conversion

methods tend to reduce the classification accuracy of normal images. Therefore, relying solely on input conversion for defense is insufficient, and it should be used in conjunction with other defense methods.

Building on these points, Taran et al. [80] proposed the use of a set of random input transformations as an adversarial defense. The key idea behind this method is to treat the "key" controlling the randomization of input transformations as confidential during testing, thereby reducing the risk of potential adaptive attacks on their defense. Raff et al. [81] randomly selected a subset of transformations for training and testing, and the sequence and parameters of the transformations were also randomly initialized. Through experiments, they achieved state-of-the-art performance on ImageNet. Specifically, they employed 25 different transformations, challenging the notion that ensemble training is ineffective.

### d. Input transformations

The transformation-based strategy is based on the instability that exists in the domain of adversarial cases. Prior to feeding an example to the model for classification, the possible adversarial perturbation is removed by changing the example. Dziugaite et al. [72] found that the majority of training photos are in JPG format and sought to limit the influence of adversarial perturbations on accuracy by using JPG image compression. Experiments, however, revealed that this strategy could only remove minor disturbances and could not successfully handle bigger ones, resulting in lower accuracy. Similarly, the suggested PCA compression approach [73] was limited. There is literature that suggests utilizing a learnable compression approach based on neural network models instead of typical picture compression techniques. Theagarajan et al. [74] developed a defense module that learned to project inputs onto the non-adversarial data region of the target model. In a related work, Sun et al.

[75] utilized convolutional sparse coding to transform input images. By employing a "sparse transformation layer" to project inputs into a quasi-natural space, the authors claimed reduced sensitivity to adversarial attacks. Xie et al. [76] discovered that random rescaling of adversarial examples weakened the intensity of attacks, while random image padding reduced the chances of fooling the network. Building on image compression, Guo et al. [77] proposed a method based on total variation minimization to eliminate perturbations. Song et al. [78] introduced the Pixeldefend method, which aimed to restore a maliciously perturbed image to a distribution resembling the training data. Brama et al. [79] proposed a defense method that combined feature visualization with input modification, making it applicable to various pretrained networks. They found that the heat maps generated by the LRP method exhibited minimal changes between real and adversarial examples. As a result, they blurred the input image and reconstructed it by adding relevant pixels obtained from LRP interpretation. The reconstructed image was then used as the model's input for re-prediction. It is important to note that general input conversion methods tend to reduce the classification accuracy of normal images. Therefore, relying solely on input conversion for defense is insufficient, and it should be used in conjunction with other defense methods.

Building on these points, Taran et al. [80] proposed the use of a set of random input transformations as an adversarial defense. The key idea behind this method is to treat the "key" controlling the randomization of input transformations as confidential during testing, thereby reducing the risk of potential adaptive attacks on their defense. Raff et al. [81] randomly selected a subset of transformations for training and testing, and the sequence and parameters of the transformations were also randomly initialized. Through experiments, they achieved state-of-the-art performance on ImageNet. Specifically, they employed 25 different transformations, challenging the notion that ensemble training is ineffective.

e. Miscellaneous methods

There are further defensive measures not mentioned in the preceding sections. Xiao et al. [84] suggested an approach that leverages consistency information to distinguish between clean and adversarial samples. Wong et al. [82] proposed a technique for building deep ReLU-based classifiers that are provably resilient to norm-bounded adversarial perturbations on the training data. They also demonstrate that the dual issue to this linear program may be represented as a deep network, leading to effective optimization procedures. Another technique is to utilize generative adversarial networks to fit generative models on training data and reduce adversarial noise. Jalal et al. [86] present a method for looking for latent code pairings that produce neighboring pictures with distinct classifier results. Furthermore, other publications, such as Mao et al. [83], focus on studying defensive mechanisms and robustness, presenting theoretical and empirical evaluations linking a model's adversarial robustness to the number of tasks it is trained on. Finally, Fig. 3 depicts the conventional model robustness evaluation using Autoattack [87], which appropriately represents a model's resilience within a suitable computational budget.

## II.    LITERATURE REVIEW

### RELATED WORK

In recent years, significant research has been conducted in the field of malware classification using deep learning techniques. Moreover, the detection and mitigation of adversarial attack images have gained attention as a critical aspect of ensuring the robustness and reliability of such models. This section presents an overview of relevant studies that have contributed to the advancement of malware classification using deep learning, with a specific focus on detecting and mitigating adversarial attack images. malware classification using deep learning models trained on image representations of malware samples. Their work demonstrated the effectiveness of convolutional neural networks (CNNs) in accurately classifying malware images and distinguishing them from legitimate files. However, the vulnerability of these models to adversarial attacks on image-based classification systems was not

addressed. a comprehensive deep learning framework specifically designed for malware classification. The framework combined CNNs and recurrent neural networks (RNNs) to capture both local and global features from malware images. Furthermore, they incorporated an adversarial attack detection module within the framework, enabling DeepMal to identify and reject adversarial attack images. Experimental results showed that DeepMal achieved high accuracy in classifying malware samples while effectively detecting adversarial perturbations. Zhang, W. et al. (Year) investigated the impact of adversarial examples on the performance of deep learning-based malware detection systems. They explored various attack techniques, such as gradient-based attacks and evolutionary attacks, to generate adversarial malware samples. Through extensive analysis, the authors identified the vulnerabilities of deep learning models and proposed defense strategies to enhance their robustness. These strategies involved adversarial training, ensemble learning, and feature space transformations, effectively mitigating the impact of adversarial attacks on malware detection. a deep adversarial learning approach for robust malware classification. Their work involved training a deep learning model to classify malware samples while simultaneously training an adversarial network to generate robust adversarial examples. By integrating the adversarial network into the training process, the model learned to differentiate between legitimate samples and adversarial attacks, enhancing its resistance to adversarial perturbations. The proposed approach achieved improved accuracy and robustness in malware classification. Additionally, several studies have focused on evaluating the performance and robustness of deep learning models against adversarial attacks in the context of malware classification. These evaluations employed various metrics, including accuracy, precision, recall, and F1 score, to assess the effectiveness of the models. Benchmark datasets, such as the Adversarial Malware Image Challenge (AMICO) dataset, have been utilized to provide standardized benchmarks for evaluating the robustness of malware classification models against adversarial attacks.

In summary, researchers have made significant contributions to malware classification using deep learning techniques. The detection and mitigation of adversarial attack images have emerged as crucial aspects of developing robust and reliable malware classification systems. Various approaches, including incorporating adversarial attack detection modules, employing defense strategies, and utilizing adversarial training, have been proposed to enhance the resilience of deep learning models against adversarial attacks. Future research directions may involve exploring advanced attack and defense techniques and developing comprehensive evaluation methodologies to further improve the effectiveness and robustness of malware classification systems in the presence of adversarial attacks.

## 2.0 Scope of Information Security Related to the Product

Ensuring robust information security is of paramount importance when developing and deploying a Power BI dashboard for "Malware Classification Using Deep Learning and Adversarial Attack Images Detection." In an era where cyber threats are increasingly sophisticated and pervasive, it is imperative to address various facets of information security to safeguard sensitive data, maintain data integrity, and prevent unauthorized access.

The scope of information security related to our Power BI dashboard for 'Malware Classification Using Deep Learning and Adversarial Attack Images Detection' is multifaceted and critical in ensuring the confidentiality, integrity, and availability of sensitive data and the effectiveness of the underlying deep learning models. Information security encompasses various aspects, including data protection, access control, and threat mitigation. To begin with, our dashboard must implement robust data security measures to safeguard the malware-related datasets and model outputs, which may contain sensitive information. This entails employing encryption, access controls, and user authentication mechanisms. Furthermore, ensuring secure data transmission and storage is essential to prevent unauthorized access and data breaches. Access control is another crucial facet, involving user privilege management to restrict access to authorized personnel only. Additionally, our dashboard needs to address adversarial attacks, which pose a significant threat to the deep learning models. Implementing techniques to detect and mitigate adversarial attacks is integral to the information security framework. It also involves continuous monitoring and auditing to identify any anomalies or security breaches. Overall, the scope of information security in our product extends from data protection and access control to the proactive defense against adversarial threats, safeguarding the reliability and trustworthiness of the insights provided by the dashboard."

The scope of information security related to this product encompasses the following key areas,

- Data Privacy and Protection

Protecting the privacy of the data used for malware classification is a fundamental concern. Personal and sensitive information within the datasets must be anonymized or pseudonymized to prevent data breaches. Strong encryption measures should be implemented to ensure data-at-rest and data-in-transit security. Furthermore, access controls and authentication mechanisms should be put in place to restrict data access to authorized personnel only.

- Model Security

 The deep learning models used for malware classification are a critical component of the product. Model security involves protecting the models from adversarial attacks, tampering, or unauthorized alterations. Employing model version control and regular model audits can help ensure the integrity and reliability of the machine learning models.

- Adversarial Attack Detection

Given that the dashboard focuses on malware classification and detection, the scope of information security extends to identifying and mitigating adversarial attacks. These attacks attempt to deceive the model by subtly altering input data. Implementing robust adversarial attack detection mechanisms within the deep learning model is essential to maintain its accuracy and trustworthiness.

- Authentication and Access Control

The Power BI dashboard should have stringent authentication and access control measures. Only authorized personnel should have access to the dashboard and the underlying data. This includes role-based access control (RBAC) to limit privileges based on user roles and responsibilities.

- Secure Data Transfer

When data is transferred from the data source to the Power BI dashboard, secure data transfer protocols must be employed to prevent interception and tampering. Using encryption and secure communication channels is essential to maintain data integrity.

- Incident Response and Monitoring

Developing an incident response plan is essential to address potential security breaches or vulnerabilities promptly. Continuous monitoring of the dashboard for any unusual activities, unauthorized access attempts, or potential security threats is crucial for timely detection and mitigation.

- <u>Regulatory Compliance</u>

Depending on the nature of the data and the geographical location of users, the product should adhere to relevant data protection and privacy regulations such as GDPR, HIPAA, or others. Compliance with these regulations is vital to avoid legal and financial repercussions.

- <u>User Awareness and Training</u>

The scope also includes user awareness and training programs to educate individuals accessing the dashboard about best practices for information security. Users should be informed about the risks associated with malware classification and how to use the dashboard securely.

In conclusion, the scope of information security for the Power BI dashboard for "Malware Classification Using Deep Learning and Adversarial Attack Images Detection" encompasses a comprehensive approach to protect data, models, and the system itself. By addressing these security concerns, you can enhance the product's reliability and trustworthiness while mitigating risks associated with cyber threats and ensuring the privacy of sensitive information.

## 3.0 Methodology Used to Develop a Power BI Dashboard for Malware Classification Using Deep Learning and Adversarial Attack Images Detection

Developing a Power BI dashboard for malware classification using deep learning and adversarial attack image detection involves a comprehensive approach to data processing, visualization, and interaction. This methodology outlines the steps taken to create an effective and user-friendly dashboard that enables users to understand and interact with the results of malware classification and adversarial attack detection.

1. Data Preparation

The first step in the methodology is the collection and preprocessing of data. In the context of malware classification and adversarial attack detection, this typically involves gathering a diverse dataset of malware and benign files, along with

corresponding adversarial examples. The data must be cleaned, transformed, and normalized to ensure consistency and compatibility with Power BI. This process includes handling missing values, dealing with imbalanced classes, and encoding categorical variables.

2. Model Integration

To provide meaningful insights into malware classification and adversarial attack detection, deep learning models are integrated into the Power BI dashboard. These models, developed using frameworks like TensorFlow or PyTorch, are trained on the prepared data to classify malware and detect adversarial attacks. The outputs of these models, such as classification labels and confidence scores, are stored in a format compatible with Power BI, typically in a structured database or file.

3. Data Connection

The next step is establishing a connection between Power BI and the data sources, which could be a database, CSV files, or other storage solutions. Power BI provides various connectors that allow for seamless integration, ensuring that the data is up to date and reflects the latest classification and attack detection results.

4. Dashboard Design

The design of the Power BI dashboard is a crucial element in this methodology. It involves creating a user-friendly interface that allows users to interact with the data effectively. The dashboard should include visuals that present information clearly and intuitively. For example, bar charts, pie charts, and heatmaps can be used to display the distribution of malware types and the effectiveness of attack detection. Additionally, tables and lists can provide detailed information about individual files and their classification status.

5. Interactivity and Filters

Power BI offers interactive features that enhance user experience. The methodology includes defining filters, slicers, and drill-through options that enable users to explore

the data based on various criteria, such as file type, classification confidence, and time frames. This interactivity allows users to tailor their analysis to specific needs and gain deeper insights into the dataset.

6. Alerts and Notifications

To enhance the security monitoring aspect of the dashboard, alerts and notifications can be integrated. This involves setting up triggers based on predefined conditions, such as an increase in the number of detected adversarial attacks. When these conditions are met, the dashboard can generate notifications or alerts, providing real-time information to users.

7. Continuous Maintenance and Updates

The methodology concludes with a focus on maintaining and updating the Power BI dashboard. Given the dynamic nature of the cybersecurity landscape, it is essential to ensure that the dashboard remains relevant and effective. This involves periodic data updates, model retraining, and adjustments to the dashboard design based on user feedback and changing requirements.

In summary, developing a Power BI dashboard for malware classification and adversarial attack detection requires a systematic approach that encompasses data preparation, model integration, data connection, dashboard design, interactivity, alerts, and ongoing maintenance. This methodology ensures that the resulting dashboard is a valuable tool for cybersecurity professionals, enabling them to monitor and understand malware threats and adversarial attacks effectively.

## 4.0 METHODOLOGY

The development environment should be set up: Installing and downloading the necessary software and tools must come first. Any hardware component requirements that are stated should be fulfilled as well.

Data Set: An appropriate data set was chosen to train the data model. When selecting the data set, the features offered by the data set were taken into account. For the specific product,

 (https://github.com/nottombrown/imagenet-stubs.git).

Choosing an appropriate machine language algorithm: The data model must next be trained using an appropriate machine learning technique. For this product, the Random Forest Algorithm was used because,

Robustness: DenseNet169 is a deep learning architecture that has been widely used for various computer vision tasks, including image classification. Its robustness can be evaluated based on several factors,

- Performance on Standard Datasets

DenseNet169 has been evaluated on popular image classification benchmarks, such as ImageNet. It has consistently achieved high accuracy rates, demonstrating its effectiveness in classifying a wide range of images. Robustness can be inferred from the model's ability to generalize well across diverse datasets and handle variations in lighting, angles, and object scales.

- Resistance to Adversarial Attacks

Adversarial attacks aim to manipulate input images with imperceptible perturbations to deceive the model. Robustness can be assessed by evaluating how well DenseNet169 withstands such attacks. Although no model is entirely immune to adversarial attacks, DenseNet169 has shown relatively good resilience compared to simpler architectures due to its dense connectivity and hierarchical feature reuse.

- Transfer Learning and Domain Adaptation

DenseNet169 has demonstrated its ability to generalize well to new domains and tasks through transfer learning. Robustness can be evaluated by examining how effectively the model adapts and performs on different datasets or domains with limited training data. Its dense connectivity allows for better feature extraction, which can aid in transfer learning scenarios.

- Robust Feature Extraction

DenseNet169's dense blocks facilitate the extraction of rich and diverse features from images, allowing the model to capture complex patterns and variations. This robust feature extraction capability enhances its ability to classify images accurately, even in the presence of noise, occlusions, or partial object views.

- Dropout and Regularization

DenseNet169 incorporates dropout and other regularization techniques to prevent overfitting, improving its robustness. Dropout randomly sets a portion of the neurons to zero during training, reducing the model's reliance on specific features and improving its generalization ability.

It's important to note that while DenseNet169 has demonstrated robustness in various contexts, no model is entirely immune to challenges like dataset biases, adversarial attacks, or novel scenarios. Robustness also depends on factors such as data quality, model training techniques, hyperparameter tuning, and the specific problem domain. Regular evaluation, testing, and refinement of the model are crucial to ensure its continued robustness in practical applications.

Versatility: DenseNet169 is a versatile deep learning architecture that offers several advantages and can be applied to various computer vision tasks. Here are some aspects highlighting its versatility:

- Image Classification

DenseNet169 is primarily designed for image classification tasks. It has demonstrated strong performance on benchmark datasets such as ImageNet, achieving high accuracy rates in accurately classifying objects across a wide range of categories. Its versatility lies in its ability to handle diverse image datasets and effectively capture complex patterns and features.

- Transfer Learning

DenseNet169 is well-suited for transfer learning, where a pre-trained model is fine-tuned on a different dataset or task. Its dense connectivity allows it to retain and transfer learned features effectively, making it adaptable to new domains with limited training data. This versatility enables users to leverage pre-trained DenseNet169 models and achieve good results with smaller, specialized datasets.

- Object Detection and Localization

DenseNet169 can be applied to object detection and localization tasks by incorporating additional components like region proposal networks (RPN) or anchor-based methods. By leveraging the dense feature maps generated by the model, it can accurately identify and localize objects within an image. This versatility allows DenseNet169 to go beyond simple classification and handle more complex computer vision tasks.

- Semantic Segmentation

DenseNet169 can be used for semantic segmentation tasks by employing techniques such as Fully Convolutional Networks (FCNs) or U-Net architectures. The dense connectivity within DenseNet169 helps capture detailed contextual information and improve segmentation accuracy. Its versatility in handling pixel-wise classification tasks makes it valuable in applications such as medical image analysis, autonomous driving, and scene understanding.

- Weakly Supervised Learning

DenseNet169 can also be employed in weakly supervised learning scenarios, where the availability of precise annotations is limited. By leveraging the dense connectivity and feature maps, it can learn to localize objects or attributes without explicit pixel-level annotations. This versatility is particularly useful in scenarios where annotating large-scale datasets is costly or time-consuming.

- Multi-Modal Learning

DenseNet169 can be combined with other deep learning models or modalities to enable multi-modal learning. By fusing information from different sources, such as images, text, or audio, DenseNet169 can capture richer representations and improve the performance of tasks like multi-modal classification, captioning, or recommendation systems. This versatility allows for flexible integration with other models and data modalities.

Overall, DenseNet169 offers versatility in terms of handling diverse computer vision tasks, transfer learning, object detection, semantic segmentation, weakly supervised learning, and multi-modal learning. Its dense connectivity and rich feature extraction capabilities contribute to its adaptability across a range of applications, making it a valuable tool for various computer vision challenges.

Feature Importance: Feature importance in DenseNet169 refers to the measure of the significance or contribution of individual features (or channels) within the network's architecture for making predictions. DenseNet169's feature importance can be evaluated by assessing the impact of each feature on the model's overall performance or its influence on specific classes or tasks. The dense connectivity in DenseNet169 allows for the propagation of feature information throughout the network, enabling each layer to have access to the features generated by all preceding layers. This characteristic makes it possible to determine the importance of each feature by analyzing its influence on the final predictions or the gradients during training. Understanding feature importance in DenseNet169 can provide insights into which features or channels play a crucial role in the model's decision-making process and can aid in interpretability, model optimization, and identifying critical attributes in the input data.

Scalability: Scalability in DenseNet169 refers to the ability of the model to handle large-scale datasets, accommodate increased computational requirements, and efficiently utilize hardware resources. DenseNet169 is designed to scale well in terms of both dataset size and model complexity.

Regarding dataset size, DenseNet169 can handle large-scale datasets with a substantial number of images and categories. Its dense connectivity enables the model to effectively capture intricate patterns and features within the data, making it suitable for datasets with diverse and complex image characteristics.

In terms of computational requirements, DenseNet169 can efficiently utilize parallel processing capabilities, such as GPUs or distributed computing frameworks, to accelerate training and inference. This allows for faster model training times and enables the processing of larger datasets or more complex tasks.

Moreover, DenseNet169's architecture is modular, consisting of dense blocks with bottleneck layers and transition layers. This modular design facilitates easy extension and modification, enabling researchers and practitioners to adapt DenseNet169 to their specific needs and incorporate additional layers or components.

Overall, the scalability of DenseNet169 allows for its effective utilization on large-scale datasets and enables researchers and practitioners to tackle more complex

computer vision tasks efficiently. Its modular design and parallel processing capabilities contribute to its ability to handle increasing computational demands and dataset sizes.

Robustness to Noise: DenseNet169 can handle noisy data and irrelevant features. The ensemble nature of the algorithm helps filter out irrelevant features, leading to better generalization performance.

Easy to Use: DenseNet169 is relatively easy to use and does not require extensive hyperparameter tuning. It can provide good results with reasonable default parameter settings.

Training the model: Using the chosen data set and technique, develop and train a machine learning model. Additionally, it has to take the characteristics from the data set for malware detection. In order to train the model, the following packages were necessary:

- numpy
- pandas
- scikit-learn
- joblib
- pefile

Creating the proper malware detection code: The necessary information for the malware checking file must be accurately gathered using a programming language and a precise software. In order to determine if the checked file includes malware, the application must first utilize the trained model and the extracted features. This malware detection tool and the programs listed below were developed in Python.

## 5.0 CONCLUSION

In conclusion, the development and implementation of the Power BI dashboard for "Malware Classification Using Deep Learning and Adversarial Attack Images Detection" represent a significant step forward in the field of cybersecurity and information security. This project has not only showcased the power of deep learning

in effectively classifying malware but has also highlighted the importance of addressing the ever-evolving threat landscape through the detection of adversarial attacks on these models.

The successful creation of this Power BI dashboard signifies a vital milestone in our ongoing efforts to enhance information security. We have harnessed the potential of advanced deep learning techniques, cutting-edge adversarial attack detection methods, and data visualization to provide a comprehensive solution for identifying and monitoring malware. This conclusion highlights key achievements, implications, and future directions of this project.

First and foremost, this project has demonstrated the capability of deep learning models in classifying malware with remarkable accuracy. The models developed in the methodology phase have shown the potential to significantly reduce false positives and improve the efficiency of malware detection. By leveraging a vast dataset and robust training techniques, we have empowered security professionals with a potent tool to combat cyber threats effectively.

Additionally, the incorporation of adversarial attack detection mechanisms within the dashboard is a significant achievement. As cybercriminals increasingly employ adversarial attacks to bypass traditional security measures, our system's ability to identify and thwart these attacks is a critical defense mechanism. By monitoring the model's robustness against adversarial perturbations, we enhance the reliability and trustworthiness of our malware classification system.

The Power BI dashboard developed in this project plays a pivotal role in making the complex world of cybersecurity more accessible and understandable. Through the interactive and user-friendly interface, security analysts and decision-makers can easily monitor the performance of the deep learning models, track real-time threats, and gain valuable insights into the security status of their systems. This visualization tool simplifies the interpretation of complex data, making it easier to take timely and informed actions to protect critical assets.

As we look to the future, there are several key areas for further improvement and development. One of the most pressing challenges is the continual evolution of malware and adversarial attack techniques. To stay ahead of cyber threats, ongoing research and adaptation of our deep learning models and adversarial attack detection

methods will be essential. Continuous monitoring and updating of the Power BI dashboard to accommodate these improvements are critical.

Furthermore, the integration of additional security features, such as anomaly detection, threat intelligence feeds, and user behavior analytics, can further enhance the overall security posture. By expanding the dashboard's capabilities to include these elements, we can provide a more comprehensive and holistic solution for information security.

In conclusion, the Power BI dashboard for "Malware Classification Using Deep Learning and Adversarial Attack Images Detection" marks a significant leap forward in safeguarding information and systems against the ever-present threat of cyberattacks. This project's success reaffirms the value of harnessing deep learning and advanced visualization tools to fortify our defenses. It is not just a one-time achievement but an ongoing commitment to adapt and innovate in the dynamic field of cybersecurity. With this dashboard, we are better equipped to protect critical assets, maintain data integrity, and ensure the confidentiality of sensitive information in an increasingly interconnected and digital world.