

Uncovering Adversarial Vulnerabilities in Deep Learning: A Shapley Value-Based Approach for Detailed Pixel Analysis

Arooj Arif

Abstract—Adversarial attacks in the deep learning domain significantly threaten model transparency and robustness. Adversarial perturbations are unnoticeable for humans. Hence, it is necessary to develop methods of counteracting these attacks. In this research, we present a new way to use SHAP (SHapley Additive exPlanations) to analyze image pixels at a finer level. We thoroughly investigate four adversarial attacks, carefully studying their effects on model vulnerability by analyzing their impact at different epsilon values. Our approach focuses on identifying crucial pixels in image-based model, uncovering vulnerabilities, and enhancing model robustness through SHAP analysis. This extensive study provides valuable insights for the development of AI systems that are more secure and transparent. These findings have important implications in critical healthcare and autonomous driving areas. Our research solves existing problems with AI security and establishes new benchmarks for how AI systems can combine robustness with interpretability.

I. INTRODUCTION

The field of deep learning is currently experiencing a golden era and has become pivotal for modern computational applications. Deep learning is a type of machine learning that uses multiple layers to process data and build computational models. It has introduced various algorithms, including generative adversarial networks, convolutional neural networks, and model transfers. These algorithms have revolutionized the way we process information. Deep learning has shown significant progress in various domains, such as visual, audio, and text processing, social network analysis, and natural language processing. It has also successfully tackled challenges in machine learning, such as unsupervised and online learning. Deep learning's ability to handle massive and complex datasets has made it a critical tool for big data analysis. Compared to traditional machine learning approaches, its state-of-the-art performance has led to its widespread adoption in fields like image processing, computer vision, speech recognition, machine translation, medical imaging, and many others [1], [2]. In the present scenario of widespread technological advancements and adoption, it has become essential to concentrate on the robustness of deep learning

models. These models are increasingly utilized in safety-critical and socially significant applications such as autonomous driving, face recognition, and malware detection, where their reliability and performance are of utmost importance. The deep neural networks have shown to be vulnerable to both adversarial and natural image corruptions that can substantially reduce their efficiency. By exploring and enhancing the robustness of these models, we can ensure their effectiveness and reliability in real-world situations [3], [4].

Adversarial attacks present a unique challenge for deep learning models. They exploit subtle weaknesses, causing misclassification of examples with imperceptible changes. Such attacks can have a severe impact on the reliability and safety of deep learning models, especially in safety-critical applications like autonomous driving. Adversarial examples can lead to incorrect decisions, potentially resulting in severe consequences. Therefore, it is crucial to analyze the robustness and reliability of deep learning models to ensure their safety in real-world deployments [5], [6]. Explainable Artificial Intelligence (XAI) is a critical component in addressing the challenges of deep learning. It helps to enhance the understanding and trustworthiness of AI by providing interpretable and human-understandable explanations of AI decisions. The techniques, tools, and algorithms used in XAI generate explanations that help build trustworthy and interpretable deep learning models. These explanations improve trust by providing insights into the model's decision-making process, addressing challenges of trust, transparency, bias understanding, and fairness, and promoting a more robust and impartial decision-making process [7], [8], [9].

Despite recent advancements in deep learning models, research has identified gaps and shortcomings in making these models robust against adversarial attacks. It has been observed that there is an overly disproportionate focus on adversarial machine learning compared to non-adversarial robustness. Additionally, there is a significant gap in model performance when

faced with naturally-induced image corruptions or alterations, which can result in performance degradation similar to that seen in adversarial conditions. This vulnerability to natural image corruptions suggests that understanding model performance on natural data should be prioritized before focusing on resilience to adversarial attack scenarios [3], [10], [11].

We conducted a research to explore the challenges posed by adversarial attacks on deep learning models. Our primary focus was on understanding and quantifying the impacts of different pixels. To achieve this, we used SHAP (SHapley Additive exPlanations) within the realm of Explainable Artificial Intelligence (XAI) to unravel the complex dynamics of how deep learning models respond to adversarial manipulations. Our goal was not only to identify vulnerabilities but also to enhance the robustness and transparency of these models. The stakes are high in critical applications such as healthcare and autonomous systems. Therefore, our research strives to fortify these models against adversarial threats while simultaneously improving their interpretability and trustworthiness. Through this work, we aim to bridge the existing gap between robustness and transparency in AI, offering novel insights and methodologies that could significantly advance the field.

II. CONTRIBUTIONS OF THE PAPER

This paper introduces several significant contributions to the field of AI, particularly in enhancing the robustness and transparency of deep learning models against adversarial threats:

- **Novel Integration of SHAP-based XAI in Adversarial Analysis:** Our work pioneers the application of SHAP-based critical pixel analysis, offering a new perspective in understanding and mitigating these threats at a pixel level.
- **Comprehensive Evaluation of Model Behavior:** By employing the MNIST dataset across a range of adversarial intensities, our research provides a detailed assessment of model vulnerabilities, crucial for developing more resilient AI systems.
- **Insightful Analysis through Advanced Visualization:** Utilizing UMAP visualizations, our study reveals intricate patterns and impacts of adversarial attacks, enhancing the interpretability of complex model behaviors in a user-friendly manner.
- **In-depth Contrastive Analysis of SHAP Values:** Our research conducts a thorough comparison of SHAP values between normal and adversarial examples, shedding light on the subtle ways adversarial attacks influence model decision-making.
- **Statistical Validation of Model Vulnerabilities:** We introduce a rigorous statistical approach to

validate the significance of identified vulnerabilities, thereby strengthening the reliability of our findings and setting a new standard in adversarial robustness research.

III. RELATED WORK

In the evolving landscape of Explainable Artificial Intelligence (XAI) and its application in mitigating adversarial attacks on AI models, a collective review of prominent research reveals a tapestry of distinct methodologies, objectives, and contributions. These studies not only underscore the diversity in approaches but also highlight the dynamic nature of the field, each carving out its unique niche.

The study "Explainable AI for Inspecting Adversarial Attacks on DNNs" offers a comprehensive analysis of the impact of adversarial attacks on deep neural networks, focusing on enhancing the overall interpretability of models. It contrasts with "Attack-agnostic Adversarial Detection on Medical Data Using Explainable ML," which zeroes in on the medical domain, combining advanced machine learning techniques with explainability to detect adversarial threats. While both employ XAI, their domain-specific applications showcase the versatility of explainable techniques in different contexts.

"PatchAttack: A Black-box Texture-based Attack with Reinforcement Learning," on the other hand, diverges by exploring the generation of texture-based adversarial attacks. This study delves into a more focused aspect of adversarial methodology, utilizing reinforcement learning for attack generation, thus contributing a new dimension to understanding and crafting adversarial strategies.

"Deep Learning under Privileged Information Using Heteroscedastic Dropout" ventures into another specialized domain, addressing the challenges of handling privileged information in deep learning models. This study's innovative application of heteroscedastic dropout highlights the ongoing efforts to adapt AI models to complex, real-world data scenarios.

Against this backdrop, your paper, "Analyzing Adversarial Vulnerabilities in Deep Learning: Unveiling Shapley Insights for Fine-Grained Pixel Analysis," emerges with a unique proposition. It delves deeper into the granularity of image data, employing SHAP analysis to unravel the influence of individual pixels in adversarial contexts. This microscopic approach is pioneering, as it not only broadens the understanding of adversarial vulnerabilities but also paves the way for developing AI models that are robust and secure against nuanced adversarial tactics. Your work stands as a testament to the innovative fusion of XAI with

AI security, underscoring the importance of detailed analysis in enhancing both transparency and resilience in AI systems.

IV. PROPOSED METHODOLOGY:

The goal of our suggested model, which is shown in Figure 1, is to improve deep learning systems' resistance to adversarial assaults. Using the MNIST dataset as a training dataset, a pre-trained Convolutional Neural Network is used to focus on extraction of correct examples. The robustness of the model is then evaluated by subjecting it to exhaustive adversarial assault simulations through a variety of methodologies. Subsequently, a comprehensive SHAP analysis is conducted to gain insights into the specific impacts of individual pixels on the model's decision-making process. In the last phase, XAI signatures are used to identify critical pixel for attack generation and detection. The implementation of this multi-stage strategy guarantees a thorough and resilient protection against adversarial threats in machine learning models.

V. STAGE 1: DATA AND MODEL PREPARATION

During the initial stage of our inquiry, we employed a convolutional neural network (CNN) that had been pre-trained on the MNIST dataset to classify digits. CNN's well-developed expertise in extracting features from picture data was essential for our investigation. One crucial preprocessing step entailed the selection of cases that were effectively classified by the model, resulting in the creation of a dataset consisting of precisely anticipated instances. The objective of this strategy was to separate the impacts of adversarial perturbations on a model that is otherwise behaving appropriately.

The dataset, after selection, exhibited a diverse class distribution: 973 instances with label 0, 1133 with label 1, 1016 with label 2, 989 with label 3, 969 with label 4, 882 with label 5, 937 with label 6, 1005 with label 7, 946 with label 8, and 984 with label 9. Although not perfectly balanced, this distribution mimics the natural frequency of digit occurrences in the actual world, providing a degree of scientific reliability to our study. The quality of our research depended on our acknowledging this little imbalance. Instead of testing the model in an artificially balanced environment, it enabled us to seriously evaluate its resilience to adversarial attacks in a real-world situation.

By focusing on correctly classified cases, we made sure that the next study was all about how adversarial attacks affect a model that is already good at what it does. The systematic methodology employed in the process of data selection established a well-defined and practical foundation for our inquiry. This allowed us to

thoroughly examine the ability of neural networks to withstand adversarial manipulations in a manner that closely mirrors real-world scenarios.

VI. STAGE 2: ADVERSARIAL ATTACK SIMULATION

By making adversarial examples, the goal of this step was to thoroughly test the model's stability. In order to learn how the neural network model reacts when exposed to intentionally misleading settings, this evaluation is crucial. The choice of epsilon numbers $\epsilon = [0.03, 0.04, 0.05, 0.1, 0.2, 0.25]$ was a key part of our method. These numbers were picked to show a range of disturbance intensities:

- **Lower Epsilon Values** ($\epsilon = 0.03, 0.04, 0.05$): Represent subtle but potentially effective perturbations, testing the model's sensitivity to minimal adversarial modifications.
- **Moderate Epsilon Value** ($\epsilon = 0.1$): Reports the model's performance in a moderately strong attack and does so in a fair way.
- **Higher Epsilon Values** ($\epsilon = 0.2, 0.25$): Check to see how well the model can handle more direct and damaging attempts to change it.

Choosing these values let us fully study how the model behaved under different levels of adversarial intensity, which is important for figuring out how robust it is overall. Utilized Foolbox, a Python library, to facilitate the generation of a wide range of adversarial attacks. Executed a series of adversarial attacks, including LinfBasicIterativeAttack (BIM), LinfFastGradientAttack (FGSM), LinfDeepFoolAttack, and LinfProjectedGradientDescentAttack (PGD). Systematically collected and analyzed the adversarial examples generated from these attacks.

A. Inclusion of the UMAP visualization:

In the analysis, we use a UMAP projection shown in Figure 2 to evaluate the performance of our neural network model. This includes its ability to handle adversarial examples generated by Projected Gradient Descent (PGD) attacks. The projection efficiently separates the ten-digit classes into distinct color-coded clusters (digits 0 to 9), indicating the model's pattern recognition capabilities. Figure 2 also reveals the strategic placement of adversarial examples. These examples are positioned within and around the digit clusters in a dispersed pattern, unlike the well-defined groupings of genuine digits. This pattern suggests that adversarial examples can deceive the model and lead to classification errors. The inclusion of Figure 2 in our analysis has a dual purpose. Firstly, it visually represents the deceptive nature of adversarial examples, highlighting how they can seamlessly blend with

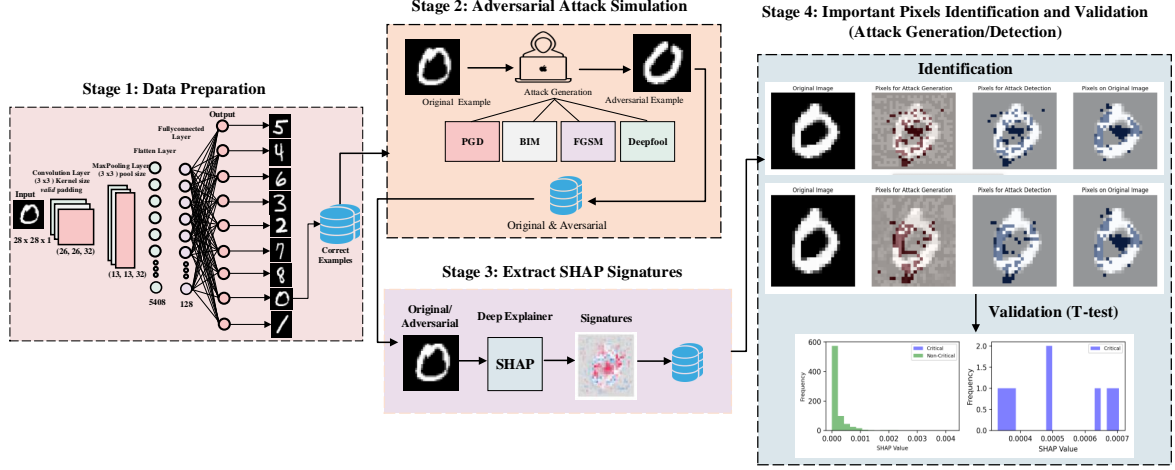


Fig. 1: Overview of the Proposed Model

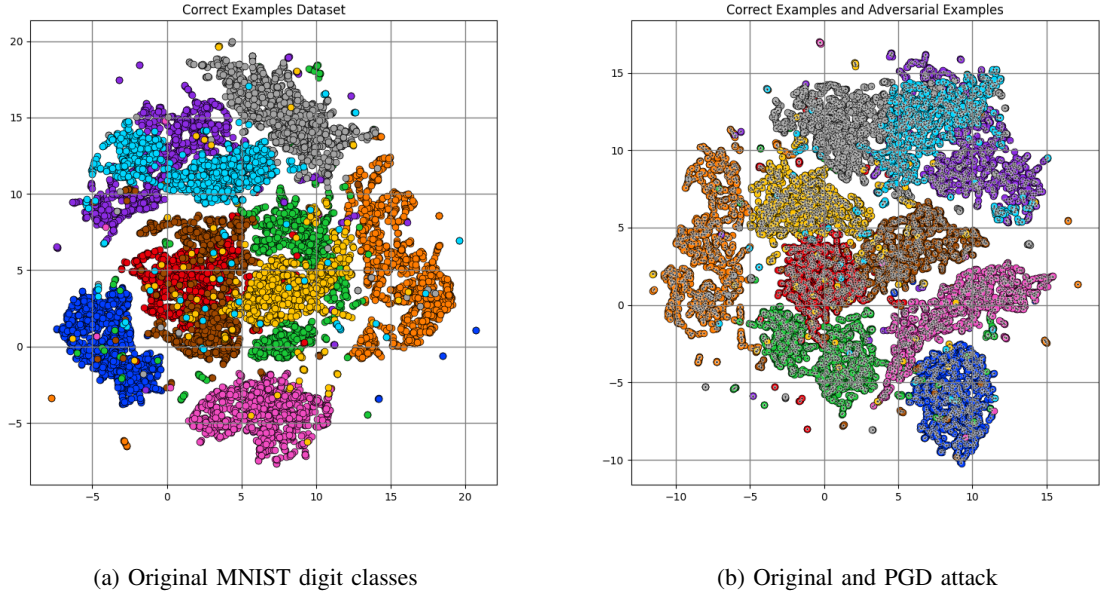


Fig. 2: UMAP projections illustrating the separation of digit classes within the MNIST dataset and the integration of adversarial examples. The left panel (2a) shows the natural clustering of MNIST digits, and the right panel (2b) overlays PGD adversarial examples.

genuine data and mislead the model. Secondly, it emphasizes the need to enhance our model's ability to differentiate between authentic and adversarial data points. This visual evidence supports our findings and underscores the urgency of fortifying neural network architectures against adversarial intrusions. Ultimately, this can help in advancing security protocols within deep learning applications.

VII. STAGE 3: EXTRACTION OF XAI SIGNATURES

A. Exploring Model Decision-Making

The main objective of this stage is to comprehend how individual input features, particularly pixels, affect the model's predictions. This process is crucial in order to distinguish the model's responses to both normal and adversarial inputs. Our approach involved:

- **Dataset Segmentation:** We selected subsets from our dataset, comprising normal examples and those modified by adversarial attacks, specifically focusing on an epsilon value of 0.2.
- **SHAP Deep Explainer Utilization:** We utilized SHAP Deep Explainer to generate XAI signatures, which expose the importance of each pixel in the model's predictions.
- **SHAP Value Calculation:** SHAP values were calculated for the selected subsets to measure the impact of each input feature on the model's output in normal and adversarial scenarios.

The SHAP values computed showed notable differences in influence patterns between normal and adversarial examples. This analysis clarifies how the model's decision-making process is altered under various input conditions.

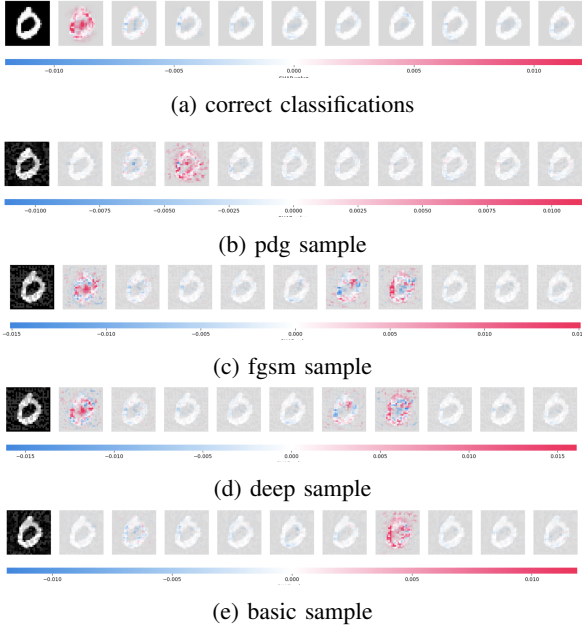


Fig. 3: SHAP heatmaps contrast the pixel influence on the model's decisions between normal and adversarial examples, underlining the adversarial impact.

VIII. STAGE 4: IDENTIFICATION AND VALIDATION OF CRITICAL PIXELS

A. SHAP Values for Critical Pixel Analysis

In Stage 4, we use the SHAP values that we extracted earlier to pinpoint and confirm the pixels that have a significant impact on the model's predictions. This step is essential for gaining a deeper understanding of the model's behavior and developing effective strategies to make it more resilient to adversarial attacks. The process involves:

Aspect	Attack Generation	Attack Detection
Objective	Identify model vulnerabilities to change predictions.	Recognize changes in interpretation from original to adversarial instances.
Focus	On influential features of adversarial examples.	On feature contribution differences between normal and adversarial examples.
Methodology	Analyze SHAP values of adversarial examples.	Compare SHAP values between normal and adversarial examples.
SHAP Values Used	From adversarial examples, indicating impactful features.	The difference in SHAP values between normal and adversarial examples.
Purpose	Craft adversarial examples exploiting vulnerabilities.	Detect and understand adversarial attacks' effects.
Insight Gained	Identifies model weaknesses and modification strategies.	Reveals decision-making changes due to attacks.

TABLE I: Key Differences Between Attack Generation and Detection Using SHAP Values

1) *Attack Generation Pixel Analysis:* We first analyze the SHAP values for adversarial image pixels. This step is crucial to identify which pixels, when altered, are most effective in generating successful adversarial attacks. By understanding the pixels that significantly influence the model's erroneous decisions, we can uncover how adversarial perturbations are guided and how they can be strategically crafted.

2) *Attack Detection Pixel Analysis:* Subsequently, we focus on detecting these adversarial manipulations by examining the differences in SHAP values between normal and adversarial images. This comparison highlights the pixels where the largest deviations occur, signaling potential areas of the image being exploited by adversarial attacks. This analysis is vital for developing robust detection mechanisms that can identify and counteract these manipulative changes.

3) *Original Image Pixel Importance:* Alongside these analyses, we also scrutinize the SHAP values of original, unmanipulated images. This examination is essential to establish a baseline of the model's interpretation of unperturbed images. Understanding the importance of pixels in the original context provides a reference point, helping us to better interpret the changes observed in the adversarial context.

Through this comprehensive methodology, we aim to provide a deep understanding of how pixel-level manipulations affect AI decision-making, both in generating and detecting adversarial attacks, as well as understanding the inherent behavior of the model under normal conditions.

B. Statistical Validation of Critical Pixels

To substantiate the distinction between critical and non-critical pixels identified in our analysis, we conducted a statistical t-test on the SHAP values. This test helps validate the significance of the differences observed in the SHAP values between critical and non-critical pixels. The implementation of a t-test ensures that the observed disparities are not due to random chance, thereby lending statistical rigor to our findings and confirming the reliability of our pixel importance assessments.

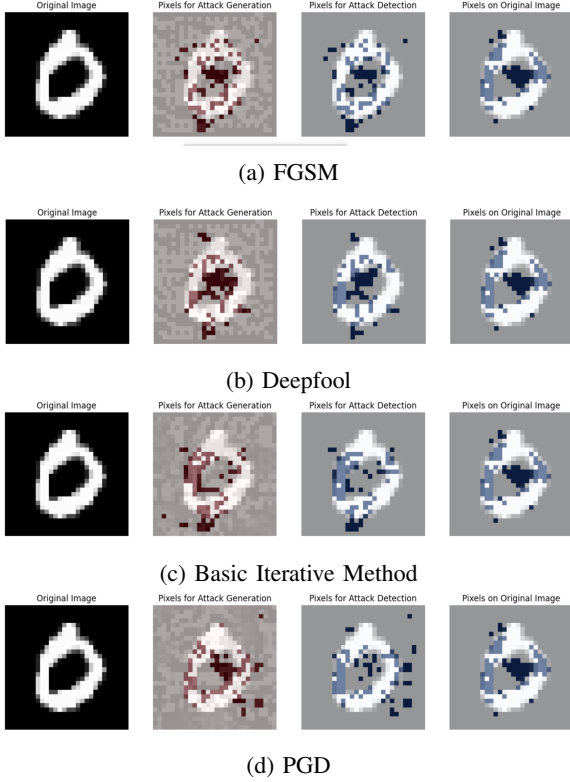


Fig. 4: Visualization of critical pixels in adversarial samples.

IX. CASE STUDIES:

This section outlines various case studies to validate the methodology.

A. Case Study 1: Evaluating Model Robustness Against Different Adversarial Attacks

The main objective of this case study is to evaluate the robustness of the MNIST classification model against different kinds of adversarial attacks such as Basic Iterative Method (BIM), Fast Gradient Sign Method (FGSM), DeepFool, and Projected Gradient Descent (PGD). The focus of the study is to compare

the performance of the MNIST model before and after being subjected to these adversarial attacks. The performance will be measured by the number of misclassifications that occur at different levels of perturbation, as measured by the epsilon values. The study aims to provide insights into the vulnerabilities of the MNIST model to different types of adversarial attacks and to determine the effectiveness of these attacks at varying intensities. Figure 5 represents the model's accuracy against various types of attacks. The graph shows how the accuracy of the model decreases with the increasing intensity of the perturbation, which is measured by epsilon. This trend of decreasing robust accuracy is observed for all types of attacks.

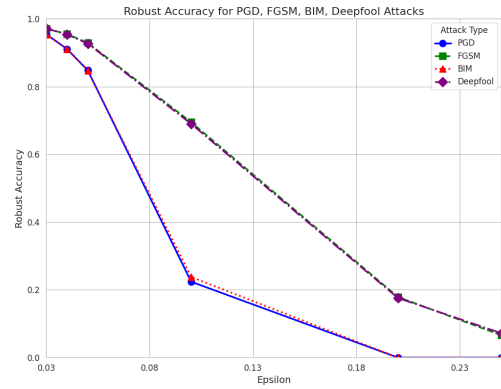


Fig. 5: Relation between attacks (PGD, FGSM, BIM, DeepFool) success rate and Perturbation size

The findings emphasize the importance of implementing advanced defensive measures to enhance the robustness of machine learning models against adversarial attacks. The observations made in this case study can help in the creation of such defenses, especially in addressing the weaknesses exposed by BIM and PGD when dealing with smaller perturbations, and the slower but eventual impact of FGSM and DeepFool when dealing with higher perturbations.

B. Critical Analysis of Class-Specific Metrics

In the field of adversarial machine learning, it's important to evaluate how well a model performs when faced with attacks. While overall metrics like robust accuracy and misclassification rates give a general idea of a model's performance, it's important to also examine class-specific metrics like F1-Score, Precision, and Recall. These metrics help to identify specific vulnerabilities in the model, especially in a security context where false positives and false negatives can have different implications. To get a more detailed

SHAP Value Proximity to Zero	Interpretation
Close to Zero	<ul style="list-style-type: none"> Feature has minimal impact on the model's prediction. Changes in the feature value do not significantly affect the prediction. Feature is considered neutral or weak in terms of prediction influence.
Far from Zero (Positive or Negative)	<ul style="list-style-type: none"> Feature has a substantial impact on the model's prediction. Variations in the feature value lead to significant changes in the prediction. Feature is a strong influencer of the predicted outcome. The magnitude of the SHAP value indicates the strength of the feature's influence.

TABLE II: Interpretations of SHAP Values Based on Proximity to Zero

TABLE III: Summary of Total Misclassifications for Various Adversarial Attacks

Attack Type	Epsilon					
	0.03	0.04	0.05	0.10	0.20	0.25
BIM	451	881	1508	7490	9834	9834
FGSM	284	440	702	2998	8070	9186
Deepfool	289	449	713	3045	8104	9130
PGD	451	881	1483	7616	9834	9834

understanding of the model's performance, we've compiled a table Table. IV that shows the F1-Score, Precision, and Recall for each class under different types of attacks and perturbations. This table not only demonstrates the model's overall resilience but also highlights specific weaknesses in certain classes that may be overlooked in more general analyses.

As an example, high precision rate during an attack indicates that when a model predicts a specific class, it is likely accurate, although it may miss out on identifying all true instances (resulting in lower recall). On the other hand, high recall with low precision suggests a model that is prone to false alarms by classifying non-members as belonging to the targeted class. The F1-Score metric balances precision and recall to provide a single measure.

Upon careful analysis of the table, it has been observed that certain classes are more susceptible to specific types of attacks. This highlights the necessity of implementing targeted defensive strategies for these classes. For instance, classes that exhibit a significant drop in F1-Score when attacked with FGSM are particularly vulnerable to the perturbations caused by this method. This knowledge can inform the creation of tailored training data or the implementation of defense mechanisms specific to these classes. Furthermore, the

variation in robustness between classes can guide the allocation of defense resources. Classes with lower robustness metrics represent weak points in the model's defenses and could benefit from focused defense measures. In conclusion, the comprehensive evaluation of the model's robustness provided by the detailed metrics, robust accuracy trends, and misclassification insights in the table is essential in developing effective defenses against adversarial attacks. This multifaceted assessment ensures that the model remains accurate and trustworthy in the face of ever-evolving adversarial challenges.

C. Pixel-Level Analysis for Enhancing AI Security Against Adversarial Attacks

1) *Analysis of Individual Pixels:* We first applied the SHAP values analysis to individual examples, focusing on the difference between normal and adversarial images. By examining a specific instance (index 3 in our dataset), we gained detailed insights into how individual pixels influenced the AI model's decision-making process. The visualization of these critical pixels (Figure 6a) reveals that certain pixels have a more pronounced impact on the model's response to adversarial attacks. This individual analysis allows us to pinpoint exact vulnerabilities in the AI model for a given input.

2) *Aggregate Analysis across Multiple Examples:* In contrast, our aggregate analysis involved averaging SHAP values across a set of images to identify commonly influential pixels. This approach provided a broader view, highlighting general patterns and shared vulnerabilities across multiple examples. The resulting heatmaps (Figure 6b) depict a more generalized perspective, essential for understanding overarching

Epsilon	Class	Metrics											
		F1-Score	Precision	Recall	F1-Score	Precision	Recall	F1-Score	Precision	Recall	F1-Score	Precision	Recall
		BIM	Deepfool	FGSM	PGD	BIM	Deepfool	FGSM	PGD	BIM	Deepfool	FGSM	PGD
0.03	0	0.9842	0.9872	0.9872	0.9842	0.9777	0.9817	0.9817	0.9777	0.9908	0.9928	0.9928	0.9908
0.03	1	0.9690	0.9788	0.9792	0.9664	0.9454	0.9601	0.9609	0.9444	0.9938	0.9982	0.9982	0.9894
0.03	2	0.9446	0.9675	0.9669	0.9442	0.9405	0.9694	0.9693	0.9396	0.9488	0.9656	0.9646	0.9488
0.03	3	0.9612	0.9778	0.9783	0.9627	0.9588	0.9749	0.9749	0.9589	0.9636	0.9808	0.9818	0.9666
0.03	4	0.9558	0.9717	0.9722	0.9573	0.9519	0.9692	0.9712	0.9548	0.9598	0.9742	0.9732	0.9598
0.03	5	0.9557	0.9715	0.9726	0.9544	0.9349	0.9571	0.9582	0.9357	0.9773	0.9864	0.9875	0.9739
0.03	6	0.9693	0.9796	0.9801	0.9693	0.9793	0.9860	0.9870	0.9783	0.9594	0.9733	0.9733	0.9605
0.03	7	0.9455	0.9682	0.9682	0.9459	0.9408	0.9682	0.9663	0.9517	0.9413	0.9682	0.9701	0.9403
0.03	8	0.9233	0.9464	0.9480	0.9234	0.9731	0.9807	0.9830	0.9720	0.8784	0.9144	0.9154	0.8795
0.03	9	0.9289	0.9550	0.9560	0.9301	0.9351	0.9619	0.9619	0.9335	0.9227	0.9482	0.9502	0.9268
0.04	0	0.9708	0.9826	0.9821	0.9718	0.9664	0.9766	0.9757	0.9693	0.9753	0.9887	0.9887	0.9743
0.04	1	0.9432	0.9708	0.9716	0.9367	0.9071	0.9464	0.9479	0.9026	0.9823	0.9965	0.9965	0.9735
0.04	2	0.8994	0.9461	0.9471	0.8980	0.8840	0.9424	0.9434	0.8822	0.9154	0.9498	0.9508	0.9144
0.04	3	0.9192	0.9596	0.9611	0.9185	0.9124	0.9577	0.9597	0.9081	0.9262	0.9616	0.9626	0.9292
0.04	4	0.9098	0.9568	0.9579	0.9116	0.8984	0.9529	0.9539	0.9028	0.9216	0.9607	0.9618	0.9205
0.04	5	0.8995	0.9546	0.9546	0.9010	0.8550	0.9320	0.9320	0.8614	0.9490	0.9785	0.9785	0.9444
0.04	6	0.9373	0.9698	0.9703	0.9402	0.9387	0.9804	0.9814	0.9579	0.9167	0.9594	0.9594	0.9232
0.04	7	0.9073	0.9512	0.9513	0.9042	0.9141	0.9521	0.9494	0.9162	0.9005	0.9502	0.9532	0.8925
0.04	8	0.8493	0.9163	0.9202	0.8527	0.9491	0.9704	0.9751	0.9447	0.7685	0.8679	0.8710	0.7770
0.04	9	0.8548	0.9304	0.9314	0.8586	0.8729	0.9371	0.9360	0.8721	0.8374	0.9238	0.9238	0.8455
0.05	0	0.9532	0.9750	0.9750	0.9559	0.9546	0.9686	0.9686	0.9549	0.9517	0.9815	0.9815	0.9568
0.05	1	0.9095	0.9534	0.9542	0.9031	0.8583	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.05	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.05	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.05	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.05	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.05	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.05	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.05	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.05	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.9069	0.9089	0.7991	0.7266	0.8618	0.8618	0.7449
0.1	0	0.9609	0.9876	0.9888	0.9601	0.9729	0.9932	0.9908	0.9781	0.9905	0.9928	0.9928	0.9830
0.1	1	0.9045	0.9518	0.9528	0.9033	0.8586	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.1	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.1	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.1	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.1	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.1	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.1	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.1	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.1	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.9069	0.9089	0.7991	0.7266	0.8618	0.8618	0.7449
0.2	0	0.9609	0.9876	0.9888	0.9601	0.9729	0.9932	0.9908	0.9781	0.9905	0.9928	0.9928	0.9830
0.2	1	0.9045	0.9518	0.9528	0.9033	0.8586	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.2	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.2	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.2	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.2	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.2	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.2	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.2	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.2	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.9069	0.9089	0.7991	0.7266	0.8618	0.8618	0.7449
0.3	0	0.9609	0.9876	0.9888	0.9601	0.9729	0.9932	0.9908	0.9781	0.9905	0.9928	0.9928	0.9830
0.3	1	0.9045	0.9518	0.9528	0.9033	0.8586	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.3	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.3	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.3	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.3	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.3	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.3	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.3	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.3	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.9069	0.9089	0.7991	0.7266	0.8618	0.8618	0.7449
0.4	0	0.9609	0.9876	0.9888	0.9601	0.9729	0.9932	0.9908	0.9781	0.9905	0.9928	0.9928	0.9830
0.4	1	0.9045	0.9518	0.9528	0.9033	0.8586	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.4	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.4	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.4	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.4	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.4	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.4	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.4	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.4	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.9069	0.9089	0.7991	0.7266	0.8618	0.8618	0.7449
0.5	0	0.9609	0.9876	0.9888	0.9601	0.9729	0.9932	0.9908	0.9781	0.9905	0.9928	0.9928	0.9830
0.5	1	0.9045	0.9518	0.9528	0.9033	0.8586	0.9169	0.9177	0.8538	0.9673	0.9929	0.9938	0.9585
0.5	2	0.8318	0.9189	0.9198	0.8318	0.8061	0.9118	0.9136	0.8105	0.8593	0.9262	0.9262	0.8543
0.5	3	0.8496	0.9311	0.9321	0.8532	0.8215	0.9260	0.9279	0.8292	0.8797	0.9363	0.9363	0.8787
0.5	4	0.8547	0.9300	0.9325	0.8589	0.8390	0.9143	0.9180	0.8481	0.8710	0.9463	0.9474	0.8700
0.5	5	0.8297	0.9208	0.9214	0.8325	0.7735	0.8825	0.8827	0.7838	0.8946	0.9626	0.9637	0.8878
0.5	6	0.8963	0.9533	0.9555	0.9047	0.9361	0.9701	0.9724	0.9402	0.8662	0.9370	0.9392	0.8719
0.5	7	0.8319	0.9297	0.9294	0.8313	0.8443	0.9320	0.9285	0.8494	0.8199	0.9273	0.9303	0.8139
0.5	8	0.7301	0.8671	0.8700	0.7372	0.8955	0.9566	0.9592	0.8739	0.6163	0.7928	0.7960	0.6374
0.5	9	0.7467	0.8838	0.8847	0.7568	0.7880	0.90						

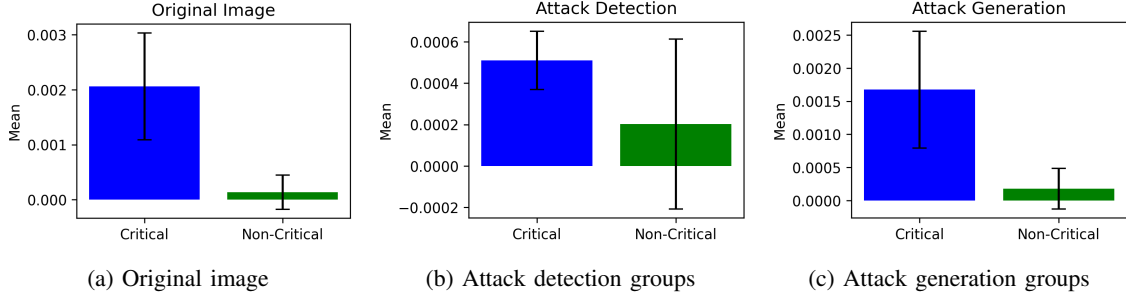


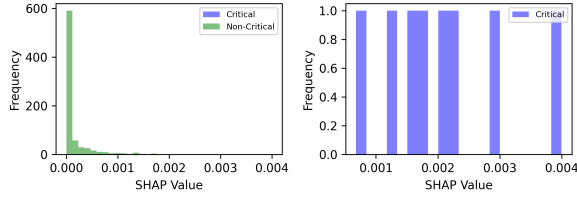
Fig. 7: The bar chart illustrates a comparison between the mean SHAP values for 'Critical' and 'Non-Critical' groups. Error bars represent standard deviations. This graph is essential for assessing the significance of the differences between these groups.

bilities of deep learning models at the pixel level. By analyzing these attacks in detail, the researchers gained insights into their intricacies. Afterwards, the power of SHAP (SHapley Additive exPlanations) analysis was used to extract Shapley values, providing a deep understanding of how each pixel contributes to model predictions. With these Shapley signatures, critical pixels within the model's decision-making process were identified. This sequential approach, starting with attack simulations, followed by Shapley analysis, and culminating in identifying critical pixels, has enriched our comprehension of adversarial threats and laid the foundation for enhancing the robustness of AI models. In conclusion, this study significantly advances our knowledge in this domain. It offers a comprehensive framework for strengthening AI models against adversarial intrusions, thereby contributing to the broader field of deep learning security.

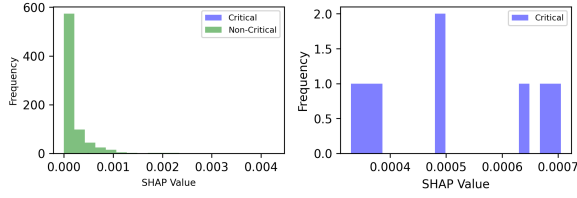
REFERENCES

- [1] S. Saad, Y. Yilin, T. Haiman, T. Yudong, P. R. Maria, S. Mei-Ling, C. Shu-Ching, and S. Iyengar. "A survey on deep learning: A. Sundaraja, "ACM Computing Surveys CSUR 51 no," 2018.
- [2] Z. Md, M. T. Tarek, Y. Chris, W. Stefan, S. Paheding, S. N. Mst, H. Mahmudul, C. V. E. Brian,
- [3] S. Numair, S. Ilya, and U. Mathias, "A systematic review of robustness in deep learning for computer vision Mind the gap," 2021.
- [4] M. Aleksandar, S. Ludwig, T. Dimitris, and V. Adrian, "Towards deep learning models resistant to adversarial attacks," 2017.
- [5] H. Samuel and N. Peyman, "Opportunities and challenges in deep learning adversarial robustness A survey," 2020.
- [6] U. Muhammad, Q. Junaid, U. J. Muhammad, A.-F. Ala, T. H. Dinh, and Niyato. "Challenges and countermeasures for adversarial attacks on deep reinforcement learning Dusit, "IEEE Transactions on Artificial Intelligence 3 no," 2021.
- [7] B. Mohammed, "Peeking inside the blackbox a survey on explainable artificial intelligence XAI," 2018.
- [8] Aha. "DARPA's explainable artificial intelligence (XAI) program David, "AI magazine 40 no," 2019.
- [9] A. Tamer, E.-S. Shaker, M. Khan, M. A.-M. Jose, C. Roberto, G. Riccardo, D. S. Javier, D.-R. Natalia, and H. Francisco, "Explainable Artificial Intelligence XAI What we know and what is left to attain Trustworthy Artificial Intelligence," 2023.

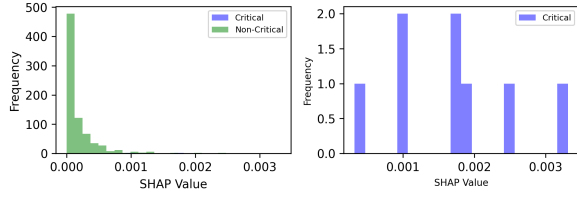
- [10] Z. Tianhang, Q. Zhan, and Liu. "Adversarial attacks and defenses in deep learning Xue, "Engineering 6 no," 2020.
- [11] E. Wei, Z. S. Quan, A. Ahoud, and Li. "Adversarial attacks on deep-learning models in natural language processing: A. survey Chenliang, "ACM Transactions on Intelligent Systems and Technology TIST 11 no," 2020.
- [12] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction," ArXiv e-prints, Feb. 2018.



(a) Original image for both 'Critical' (in blue) and 'Non-Critical' (in green) groups.



(b) 'Critical' (in blue) and 'Non-Critical' (in green) pixels for attack detection group.



(c) Adversarial image for both 'Critical' (in blue) and 'Non-Critical' (in green) pixels for attack generation groups.

Fig. 8: In the left subplot, the histogram depicts the distribution of SHAP values of Adversarial image for both 'Critical' (in blue) and 'Non-Critical' (in green) groups. The right subplot zooms in on the 'Critical' group's SHAP values. These histograms provide insights into the data distribution and its impact on statistical tests