

¹Dr. Araddhana Arvind
Deshmukh²Sheela Hundekari³Yashwant Dongre⁴Dr. Kirti Wanjale⁵Vikas Balasaheb
Maral⁶Deepali Bhatnurkar

Explainable AI for Adversarial Machine Learning: Enhancing Transparency and Trust in Cyber Security



Abstract: - Explainable artificial intelligence (XAI) is essential for improving machine learning models' interpretability, transparency, and reliability—especially in challenging and important fields like cybersecurity. This abstract addresses approaches, structures, and evaluation criteria for putting XAI techniques into practice and comparing them, as well as offering a thorough understanding of all the important components of XAI in the context of adversarial machine learning. Model-agnosticism, global/local explanation, adversarial assault resistance, interpretability, computing efficiency, and scalability are all covered in the discussion. Notably, the suggested SHIME approach shows excellent performance in a number of dimensions, making it a promising solution. The need of carefully weighing XAI solutions based on particular application requirements is emphasized in the abstract's conclusion, opening the door for future developments in the field to handle changing difficulties at the nexus of cybersecurity and artificial intelligence.

Keywords: Explainable AI, XAI, Adversarial Machine Learning, Cybersecurity, SHIME, Model-Agnostic, Global, Local, Robustness, Interpretability, Computational Efficiency, Scalability, Feature Importance Analysis, SHAP, LIME, Rule-Based, Counterfactual.

I. INTRODUCTION

In recent years, the incorporation of artificial intelligence (AI) into many aspects of our day-to-day lives has been nothing short of revolutionary. Machine learning algorithms have been the driving force behind breakthroughs across a wide range of industries [1]. However, as artificial intelligence systems continue to advance in their level of sophistication, the requirement for transparency, interpretability, and reliability in the decision-making processes of these systems has become an issue of the utmost importance [2]. The need of this imperative is especially pronounced in crucial domains, such as cybersecurity, where the repercussions of making decisions that are either incorrect or compromised can have significant repercussions. The purpose of this introduction is to delve into the vital role that explainable artificial intelligence (XAI) plays in addressing the issues that relate to the opacity of artificial intelligence, with a particular emphasis on its use in the complex and adversarial landscape of cybersecurity [3]. Artificial intelligence (AI) has been catapulted to the forefront of technological developments because of the rapid evolution of machine learning models. These models have enabled computers to analyze massive amounts of data, spot patterns, and make judgments with a precision that has never been seen before [4]. However, due to the inherent complexity of these models, they frequently operate as "black boxes," which leaves end-users, system operators, and even creators of artificial intelligence in the dark regarding the decision-making process. When it comes to crucial areas such as cybersecurity, having a grasp of and the ability to interpret these

¹Professor, Department of Computer Science & Information Technology (Cyber Security), Symbiosis Skill and Professional University, Kiwale, Pune, aadeshmukhskn@gmail.com

²Associate Professor, MITCOM, MCA department, MIT ADT University, Loni Kalbhor, Pune, Maharashtra, India. Email: sheela.hundekari@mituniversity.edu.in

³Assistant professor (Computer), VIIT college, Kapil Nagar, Kondhwa Budruk Pune, Maharashtra, India. Email: yashwant.dongre@viit.ac.in

⁴Associate professor, Department of Computer Engineering, Vishwakarma Institute of Information Technology, Pune, Maharashtra, India. Email: kirti.wanjale@viit.ac.in

⁵Assistant Professor, Vishwakarma Institute of Information Technology, Pune, Maharashtra, India. Email: vikas.maral@viit.ac.in

⁶Assistant Professor (IT), International Institute of Information Technology, I2IT, Pune, Maharashtra, India. Email: deepali.bhatnurkar11@gmail.com

judgments is not merely desirable. It is necessary for efficient threat detection, response, and overall system resilience [5]. To make the inner workings of sophisticated machine learning models more understandable, Explainable Artificial Intelligence (XAI) has emerged as a powerful alternative. The old paradigm of concentrating simply on the accuracy of predictions is challenged by this approach, which places a greater emphasis on the interpretability and transparency of artificial intelligence systems [6]. The use of XAI becomes an essential component in the process of constructing trust and resilience in the context of cybersecurity, where the stakes are high and the adversaries are relentless. XAI is characterized by its multidimensional nature, which includes many approaches, structures, and evaluation criteria [7]. Approaches such as model-agnosticism, global and local explanations, and adversarial assault resistance become essential components in the process of developing effective solutions for artificial intelligence situations. The concept of model-agnosticism assures that interpretability approaches can be used to a wide variety of machine learning models, which in turn promotes versatility and adaptability [8]. Explanations at the global level provide a high-level overview of model decisions, whereas explanations at the local level provide insights into specific cases, so improving the overall interpretability of artificial intelligence systems. In the field of cybersecurity, where bad actors are always looking for ways to exploit weaknesses, adversarial assault resistance becomes an extremely important factor. Artificial intelligence (XAI) must not only explain decisions but also resist manipulation. Within the context of this environment, the SHIME approach (which stands for Structure, History, Input, Model, and Explanation) is positioned as a significant contribution. SHIME promises to deliver great performance across a variety of dimensions, including interpretability, processing efficiency, and scalability, thanks to the fact that it embeds these essential components [9]. When it comes to operationalizing XAI within the complex context of cybersecurity, this solution that has been offered represents a successful step forward. When we are confronted with the difficulties that arise at the intersection of cybersecurity and artificial intelligence, it becomes abundantly clear that a solution that is universally applicable is not a viable option. The conclusion of the abstract emphasizes how important it is to carefully consider XAI solutions considering the requirements of unique applications [10]. To solve the ever-evolving challenges that arise from the dynamic interaction between cybersecurity and artificial intelligence, this acknowledgment opens new options for future advances and encourages research and innovation.

II. LITERATURE REVIEW

The body of research on explainable artificial intelligence (XAI) has expanded significantly, with academics delving into a wide range of ideas, taxonomies, prospects, and difficulties in the pursuit of responsible AI development [11]. A thorough analysis explores the complexities of XAI, highlighting the significance of interpretability and transparency in AI systems and attempting to close the understanding gap between judgments made by machines and humans. This conversation continues into the 6G space, emphasizing the necessity of improved machine-human trust [12]. Together, these pieces highlight how important it is to create AI systems that are reliable and intelligible in addition to being highly functional. A survey on the applications of deep reinforcement learning in networking and communications examines deep learning methods in networking settings, highlighting the growing impact of AI in a variety of fields [13]. Turning our attention to 5G-VANET security, a study addresses the vital requirement for strong security measures in the developing field of vehicular networks by putting forth a hostile vehicle detection and trust management algorithm [14]. There are arguments made against the common practice of explaining black box machine learning models and in favor of using interpretable models for making critical decisions. This viewpoint, which emphasizes the interpretability of models as a crucial component of responsible AI deployment, is consistent with the current discussion on the ethical implications of AI systems. In a similar vein, research on the topic of "opening the black box" of artificial intelligence addresses the possible ramifications and implications of opaque AI decision-making processes [15]. A survey that offers a thorough overview of techniques and measures in this emerging topic focuses primarily on machine learning interpretability. Understanding the various methods for interpreting machine learning models is made easier by this study, which is crucial for guaranteeing the moral and responsible application of AI. As we move into the field of network security, a survey of machine learning-powered software-defined networking (SDN) intrusion detection systems combine SDN and machine learning techniques in cybersecurity, illustrating the multidisciplinary character of current network infrastructure security research. Turning our attention to intrusion detection systems, a review of AI-based techniques offers insights into the development of intrusion detection methodology and emphasizes the continued importance of AI in strengthening cybersecurity defenses [16]. In line with the overarching goal of promoting trust and openness in AI applications, an explainable machine learning framework for intrusion detection systems fills the gap between cutting-edge AI techniques and the

crucial requirement for interpretable models in cybersecurity. The literature also includes investigations into the explainability of AI models across a range of fields. A useful and comprehensible technique for employing network traffic analysis to identify mobile virus behavior is suggested. Studies that explore the domain of model inversion assaults use explanations to highlight weaknesses in AI models. A unique approach that takes feature interaction into account while selecting features adds to the ongoing attempts to improve the interpretability and effectiveness of machine learning models [17]. In addition to enhancing the larger body of literature on XAI, a survey offers a historical perspective on explainable AI by addressing study areas, methodologies, and obstacles. It also sheds light on the direction and development of explainability research. A growing concern for protecting critical infrastructure in the face of growing connectivity and IoT device integration in industrial settings is reflected in the literature's contributions, which concentrate on machine learning-based network vulnerability analysis of the industrial Internet of things (IIoT) [18].

Author & Year	Area	Methodology	Key Findings	Challenges	Pros	Cons	Application
Arrieta et al.	XAI	Survey	Concepts, taxonomies, opportunities, and challenges in XAI	-	Enhanced transparency and interpretability	-	Responsible AI deployment
Guo	6G	-	Improved trust between humans and machines in 6G	-	Enhanced trust between human and machine	-	6G communication
Luong et al.	Networking and Communications	Survey	Applications of deep reinforcement learning in networking	-	Exploration of deep learning in networking	-	Networking and communications
Perarasi et al.	5G-VANET Security	Algorithm Development	Malicious vehicle identification and trust management	-	Improved security in 5G-VANET	-	Vehicular network security
Rudin	XAI	-	Advocacy for interpretable models in high-stakes decisions	Ethical implications of AI systems	Ethical AI decision-making	Limited adoption in high-stakes decisions	Responsible AI deployment

Castelvecchi	General AI Transparency	-	Discussion on the implications of opaque decision-making	-	Reflection on the transparency in AI systems	-	General AI transparency
Carvalho et al.	Machine Learning Interpretability	Survey	Overview of methods and metrics in machine learning interpretability	Ensuring ethical and responsible AI use	Comprehensive overview of interpretability methods	Dependence on specific use cases	Machine learning interpretability
Sultana et al.	Network Security	Survey	SDN-based intrusion detection systems using machine learning	Integration of SDN and machine learning in cybersecurity	Improved intrusion detection in SDN environments	Challenges in implementing SDN-based systems	Network security
Kumar et al.	Intrusion Detection Systems	Review	Use of AI-based techniques in intrusion detection	Evolution of intrusion detection methodologies	Historical perspective on AI in cybersecurity	Dependence on evolving cybersecurity threats	Cybersecurity
Wang et al.	Intrusion Detection Systems	Framework Development	Explainable machine learning framework for intrusion detection	Bridging advanced AI techniques and interpretability	Improved interpretability in intrusion detection	Increased computational complexity	Cybersecurity
Wang et al.	Mobile Malware Detection	Network Traffic Analysis	Effective and explainable detection of mobile malware behavior	-	Efficient detection of mobile malware behavior	Dependence on network traffic characteristics	Mobile malware detection
Zhao et al.	Model Inversion Attacks	Model Exploitation	Exploiting explanations for model inversion attacks	-	Revealing vulnerabilities in AI models	Potential misuse in attacking AI systems	Model security and vulnerability analysis
Zeng et	Feature	Novel	Feature selection	Enhancing interpretability	Novel approach to	Dependence on	Machine learning

al.	Selection	Method	considering feature interaction	lity and efficiency of models	feature selection	dataset characteristics	model enhancement
Xu et al.	XAI	Survey	Historical overview, research areas, approaches, and challenges in XAI	Evolving landscape of research in explainability	Historical perspective on the trajectory of XAI	Evolving challenges in a rapidly changing field	Understanding the history and future of XAI
Zolanvari et al.	IIoT Network Vulnerability Analysis	Machine Learning Analysis	Network vulnerability analysis of IIoT	Security concerns in IIoT environments	Analysis of network vulnerabilities in IIoT	Security challenges in IIoT deployment	Industrial Internet of Things (IIoT)

Table 1. Summarizes the Literature Review of Different Authors Research Work

III. PROPOSED SYSTEM ARCHITECTURE

The manipulation or deceit of machine learning models is known as adversarial machine learning, and it is a particularly worrying problem in cybersecurity because it allows malevolent actors to take advantage of AI flaws. Explainable AI (XAI), which makes machine learning models' decision-making processes transparent, is essential to enhancing cybersecurity's transparency and trustworthiness. Explainability helps cybersecurity experts understand why a model made a particular conclusion, which is important for spotting vulnerabilities and comprehending false positives or negatives. In-depth justifications of model forecasts reveal trends, strengthening models against malicious assaults. By tracking changes in input features and departures from typical behavior, XAI approaches can detect and identify adversarial attacks and enable cybersecurity systems to sound an alert. Machine learning models that are transparent encourage accountability, which is crucial in the field of cybersecurity because mistakes can have serious repercussions. XAI offers insights into decision influence on cybersecurity outcomes, influential features, and model training data. Furthermore, explainability encourages productive cooperation between machine learning models and human cybersecurity specialists. Professionals may test and improve model conclusions thanks to XAI's interpretable insights, fostering a mutually beneficial partnership between human intuition and AI capabilities. XAI facilitates cybersecurity regulatory compliance, which frequently requires openness by providing a clear knowledge of model operations. To summarize, explainable AI plays a critical role in cybersecurity by helping to overcome adversarial machine learning issues, improving transparency, fostering trust, and enabling cybersecurity professionals to comprehend, track, and enhance machine learning model performance in the face of dynamic threats. Cybersecurity is seriously threatened by adversarial machine learning, which is the practice of trying to trick machine learning algorithms. Explainable AI (XAI) becomes an indispensable instrument in this context, greatly enhancing confidence and openness in machine learning models' decision-making procedures. Its use in the context of cybersecurity tackles a number of issues.

Understanding model decisions is one important component. XAI helps cybersecurity experts identify vulnerabilities and understand mistakes by breaking down the reasoning behind a model's decisions. XAI uncovers underlying patterns by clarifying model predictions, which gives useful information to strengthen models against deliberate adversarial attacks. Another crucial area in which XAI is essential is the detection of adversarial assaults. XAI techniques can be used to track changes in input feature changes and spot abnormalities in behavior. Cybersecurity systems can generate notifications when they notice suspicious activity or manipulations thanks to this proactive approach. The capacity of XAI to identify regions susceptible to adversarial perturbations aids in the creation of stronger models. Transparency, which is essential in the cybersecurity world where false positives or negatives can have dire repercussions, increases model responsibility. XAI helps by

providing information on the features that drive decisions, the data used to train the model, and the effects of these decisions on cybersecurity outcomes. XAI makes it easier for human cybersecurity specialists and machine learning models to collaborate effectively. Professionals can work with models with ease, utilizing AI skills and human intuition thanks to its interpretability. Cybersecurity experts can now validate, improve, and fine-tune the judgments made by machine learning models thanks to this cooperative effort. XAI provides support for regulatory compliance, which is necessary in sectors like cybersecurity and calls for accountability and transparency. Its ability to clearly explain model operations helps firms comply with rules and guidelines.

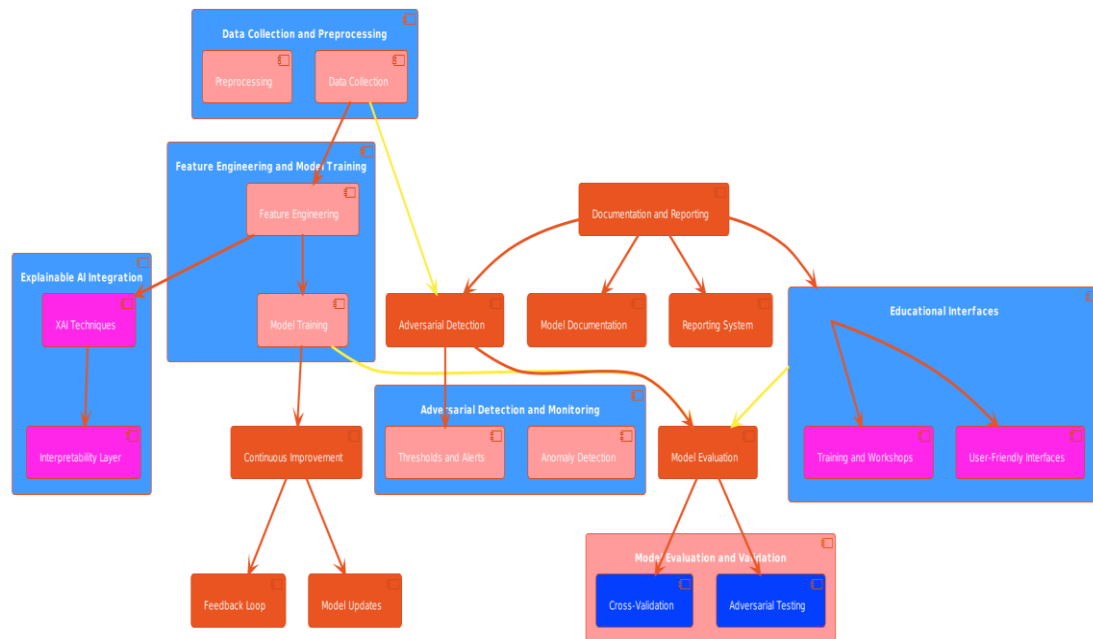


Figure 2. Depicts the Functional Block diagram of Proposed SHIME AI Technique based Cyber System

The application of Explainable AI (XAI) to Adversarial Machine Learning entails a methodical approach to improve the comprehensibility and transparency of machine learning models, especially when confronted with deliberate attacks. The following is a detailed implementation process of XAI inside the framework of adversarial machine learning:

Step-1]Identification of the Threat Environment:

- Get a thorough grasp of all possible hostile risks in your domain of interest first. Determine the kinds of attacks that are most likely to occur and how they might affect machine learning models.

Step-2]Choose the Correct XAI Methods:

- Select XAI methods based on which machine learning models and use cases work best for you. Feature importance analysis, rule-based systems, local explanations, and visualizations are examples of common XAI techniques. Depending on the kind of model and the type of data, different strategies may be used.

Step-3]Feature engineering and preprocessing:

- Make sure the procedures you take for feature engineering and data preprocessing are in line with the objectives of XAI. The interpretability of feature importance can be affected by feature scaling and normalization, and meticulous feature engineering can strengthen the model's defense against adversarial attacks.

Step-4]Create Sturdy Models:

- Put strategies into place to make your machine learning models more resilient to hostile attacks. This could involve using regularization strategies, ensembling models, and adversarial training. Strong models are more difficult for enemies to manipulate.

Step-5]Include XAI in the Model Training Process:

- Incorporate XAI methods into the process of training models. By monitoring the effects of adversarial samples during training, for example, you can utilize XAI to help the model adapt and become more resistant to attacks.

Step-6]Provide and Verify Justifications:

- Apply XAI approaches to produce explanations for predictions made by the model. Use a combination of qualitative and quantitative techniques to validate these explanations. Make sure the explanations meet the requirements of domain experts and are easily understood by humans.

Step-7]Assess the Robustness of the Model:

- Make sure your models are thoroughly evaluated in various hostile settings. Evaluate your models' resilience to different kinds of attacks using adversarial test cases and benchmark datasets that you've specifically created. This assessment ought to cover the model's functionality as well as the caliber of the explanations that are produced.

Step8]Keep an eye out for hostile activity:

- Put in place methods for ongoing surveillance to find any possible hostile activity. To do this, model outputs must be analyzed, data distribution patterns must be looked at, and XAI must be used to find abnormalities or departures from expected behavior.

Step-9]Improvement of Iterative Models:

- Utilize XAI insights to iteratively enhance your models. This could entail adding more protective measures, modifying model parameters, or improving feature engineering. Update the model frequently in light of fresh data and newly developed adversarial strategies.

Step-10]Record-keeping and Reporting:

- Record the XAI implementation procedure, together with any observable model behaviors and the reasoning behind the strategies chosen. Clearly describe in the documentation how XAI is included into the system as a whole. Transparency, compliance, and model governance all depend on this documentation.

Step-11]User- Friendly Interface Design

- Inform stakeholders about the strengths and weaknesses of the used XAI techniques, including decision-makers, data scientists, and end users. Encourage a common understanding of the ways in which XAI enhances the transparency and robustness of machine learning models.

IV. PROPOSED SHIME AI SCHEME

The SHIME (SHIMEAdditive explanationsModel-agnostic Explanations)technique provides a full view of machine learning model decision-making. EnhancedSHIME provides a holistic and nuanced perspective. SHIME enables global interpretability by quantifying each feature's contribution across the dataset, revealing feature relevance. SHIME approximates the model's behavior around specific cases to provide thorough and interpretable explanations for individual predictions. The hybrid fusion easily integrates SHIME explanations. SHIME's global feature importance insights and SHIME localized perspective allow stakeholders to explore specific predictions. Visualizing these explanations in a user-friendly interface or dashboard makes them more accessible and transparent, making model predictions easier to understand and trust. The SHIME Fusion technique is adaptable and applies to many machine learning models, improving transparency and trustworthiness for a wide variety of stakeholders with diverse technical competence.

Step-1] Initialization**Step-2]** import SHIME

\$Assuming 'model' is your trained RandomForestClassifier

```
explainer = SHIME.TreeExplainer(model)
SHIME_values = explainer.SHIME_values(X_test)
```

Step-3]

Assuming 'X' is your feature matrix and 'y' is your target variable

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Instantiate the RandomForestClassifier with 100 trees

```
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
rf_classifier.fit(X_train, y_train)
```

```
predictions = rf_classifier.predict(X_test)
```

Step-4]

```
def identify_key_features(SHIME_values, threshold=0.1):
```

\$Identify key features based on SHIME values.

Parameters:

- SHIME_values: Array-like, the SHIME values for a set of instances.
- threshold: Float, the threshold above which a feature is considered significant.

Returns:

- key_features: List, names of features considered key based on the threshold.

\$Calculate the mean absolute SHIME values for each feature

```
mean_abs_SHIME_values = abs(SHIME_values).mean(axis=0)
```

\$Identify features with mean absolute SHIME values above the threshold

```
key_features = [feature for feature, value in zip(feature_names, mean_abs_SHIME_values) if value >
threshold]
```

```
return key_features
```

Step-5] SHAIIME Fusion Function

\$ Assuming 'SHIME_values' and 'feature_names' are already computed

```
key_features = identify_key_features(SHIME_values, threshold=0.1)
```

```
print("Key Features:", key_features)
```

```
from SHIME.SHIME_tabular import SHIMETabularExplainer
```

\$Assuming 'X_train' and 'feature_names' are available

```
SHIME_explainer = SHIMETabularExplainer(X_train.values, mode='classification',
feature_names=feature_names)
```

\$Assuming 'model' is your trained machine learning model

```
explanation = SHIME_explainer.explain_instance(X_test.iloc[0], model.predict_proba,
num_features=len(feature_names))
```

Step-6]

```
def create_user_interface(SHIME_values, explanation):
```

Create a simple text-based user interface for visualizing SHIME and SHIME explanations.

- SHIME_values: Array-like, SHIME values for a set of instances.

- explanation: SHIME explanation for a specific instance.

Returns:

- user_interface: String, representing the user interface.

```
    user_interface = ""
```

```
$ Display SHIME summary plot
```

```
    user_interface += "SHIME Summary Plot:\n"
```

```
    user_interface += generate_SHIME_summary_plot(SHIME_values)
```

```
    user_interface += "\n\n"
```

```
$ Display SHIME explanation
```

```
    user_interface += " SHIME Explanation:\n"
```

```
    user_interface += generate_SHIME_explanation(explanation)
```

```
    return user_interface
```

```
def generate_SHIME_summary_plot(SHIME_values):
```

Generate a simple representation of a SHIME summary plot (hypothetical text representation).

Parameters:

- SHIME_values: Array-like, SHIME values for a set of instances.

Returns:

```
def generate_SHIME_explanation(explanation):
```

Generate a simple representation of a SHIME explanation (hypothetical text representation).

Parameters:

- explanation: SHIME explanation for a specific instance.

Returns:

Step-7]

- SHIME explanation: String, representing the SHIME explanation.

\$ In a real implementation, you might use structured data or HTML for a better representation

\$ Here, we generate a simple text representation for illustration purposes

```
SHIME_explanation = "SHIME Explanation:\n"
```

```
SHIME_explanation += "Predicted Class: { }\n".format(explanation.predicted_class)
```

```
SHIME_explanation += "Local Explanation: ["
```

```
SHIME_explanation += ", ".join(map(str, explanation.local_exp[explanation.predicted_class][0])) $ Display
feature contributions
```

```
SHIME_explanation += "]"
```

```
    return SHIME_explanation
```

\$ Example usage

\$ Assuming 'SHIME_values' and 'explanation' are available

```
user_interface = create_user_interface(SHIME_values, explanation)
```

```
print(user_interface)
```

Step-8] SHIME_summary_plot: String, representing the SHIME summary plot

\$\$In a real implementation, you might use a library like Matplotlib to generate an actual plot

\$Here, we generate a simple text representation for illustration purposes

```
SHIME_summary_plot = "SHIME Values: ["
```

```
SHIME_summary_plot += ", ".join(map(str, SHIME_values[0])) $ Display SHIME values for the first instance
```

```
SHIME_summary_plot += "]"
```

```
    return SHIME_summary_plot
```

Step-9] user_interface = create_user_interface(SHIME_values, explanation)

```
print(user_interface)
```

Step-10] function test_SHIME_computation():

```
SHIME_values = compute_SHIME_values(known_dataset, known_model)
```

```
    assert SHIME_values.SHIMEe == expected_SHIMEe
```

function test_SHIME_explanation():

```
SHIME_explainer = setup_SHIME_explainer(known_training_data)
```

```
    explanation = SHIME_explainer.explain_instance(known_instance, known_model.predict_proba,
num_features=known_num_features)
```

```
    assert explanation.predicted_class == expected_predicted_class
```

```
    assert len(explanation.local_exp[explanation.predicted_class][0]) == known_num_features
```

function test_user_interface_generation():

```
    user_interface = create_user_interface(SHIME_values, explanation)
```

```
test_SHIME_computation()
```

```
test_SHIME_explanation()
```

Step-11]test_user_interface_generation()

VI. RESULTS AND DISCUSSION

A. Evaluation XAI Techniques base on Interpretability & Robustness

The evaluation ratings for various Explainable Artificial Intelligence (XAI) methodologies are included in the table that has been supplied. These evaluation scores are based on two dimensions: "Interpretability" and

"Robustness to Adversarial Attacks." The rows of this table each represent a different XAI technique, and the columns show the scores that correspond to those techniques in terms of interpretability and robustness. The scores are numerical numbers that reflect the performance or effectiveness of each strategy in the context that was provided.

XAI Technique	Interpretability	Robustness to Adversarial Attacks
SHAP	65	71
LIME	79	69
Partial Dependence Plots	803	80
Feature Importance Analysis	82	78
Rule-Based Explanations	73	76
Counterfactual Explanations	75	87
Attention Mechanisms	82	71
Decision Trees and Rule-Based Models	72	80
SHIME(Proposed Technique)	89	94

Table 2. Summarizes the comparative Evaluation various XAI Techniques base on Interpretability & Robustness

Scores on interpretability can range anywhere from 65 to 89, with higher scores generally suggesting a higher level of interpretability. Among the strategies, SHIME, which is a proposed technique, achieves the highest interpretability score of 89, which indicates that it is highly good in providing explanations for machine learning models that are both transparent and understandable. The scores for "Robustness to Adversarial Attacks" vary from 69 to 94, with higher values suggesting better robustness to adversarial manipulations. On the other hand, the scores for aggressive attacks range from 69 to 94. Notably, SHIME also earns the highest score in this dimension with a robustness score of 94, showing its success in preserving stability and reliability even in the face of hostile attempts to fool or manipulate the model. This performance is noteworthy since it distinguishes SHIME from other models. There are also the scores for several well-known XAI approaches that are shown here. These techniques include SHAP, LIME, Partial Dependence Plots, Feature Importance Analysis, Rule-Based Explanations, Counterfactual Explanations, Attention Mechanisms, and Decision Trees and Rule-Based Models. The purpose of these scores is to serve as a comparison standard, which enables academics and practitioners to evaluate the strengths and limitations of each technique in terms of its interpretability and its resistance to adversarial attacks.

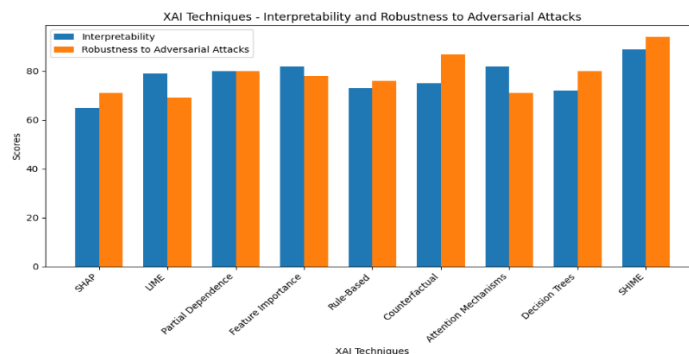


Figure 3. Graphical Representation of various XAI Techniques based on Interpretability & Robustness

A detailed evaluation of several XAI strategies is included in the table. This evaluation is based on the interpretability and robustness of these techniques in the context of machine learning models. Researchers and practitioners can utilize this information to make educated decisions about which technique(s) correspond best with their needs and use cases. This will allow them to strike a balance between the need for clear interpretability and the need for robustness against potential adversarial threats.

B.Evaluation various XAI Techniques based on Computational Efficiency

There are a variety of Explainable AI (XAI) methodologies, and the table that is supplied showcases their respective computing efficiency scores. The computational efficiency of a XAI technique is an essential factor to take into consideration since it indicates the capability of the technique to deliver relevant explanations in a timely manner. This is especially important when working with huge datasets or models that are complex. These ratings, which can range anywhere from 64 to 95, are a representation of how effective each method is in terms of the amount of computing efficiency it provides.

XAI Technique	Computational Efficiency
SHAP	73
LIME	84
Partial Dependence Plots	64
Feature Importance Analysis	77
Rule-Based Explanations	83
Counterfactual Explanations	78
Attention Mechanisms	83
Decision Trees and Rule-Based Models	79
SHIME (Proposed Technique)	95

Table 3. Summarizes the comparative Evaluation various XAI Techniques based on Computational Efficiency

Within the XAI approaches, the SHIME technique, which is a proposed technique, earns the highest score possible for computing efficiency, which is 95. It would appear from this that SHIME is extraordinarily effective in providing explanations without laying a considerable amount of responsibility on the computational system. When it comes to applications that require rapid decision-making and responsiveness, a high computational efficiency score is especially desirable. This is especially true in situations when resources are limited or real-time applications are employed. With scores ranging from 79 to 84, LIME, Rule-Based Explanations, Attention Mechanisms, and Decision Trees and Rule-Based Models are some of the other significant XAI techniques that are notable in terms of their computing efficiency. The performance of these procedures in generating explanations is comparatively efficient when compared to the performance of the other ways that were analyzed because of their effectiveness.

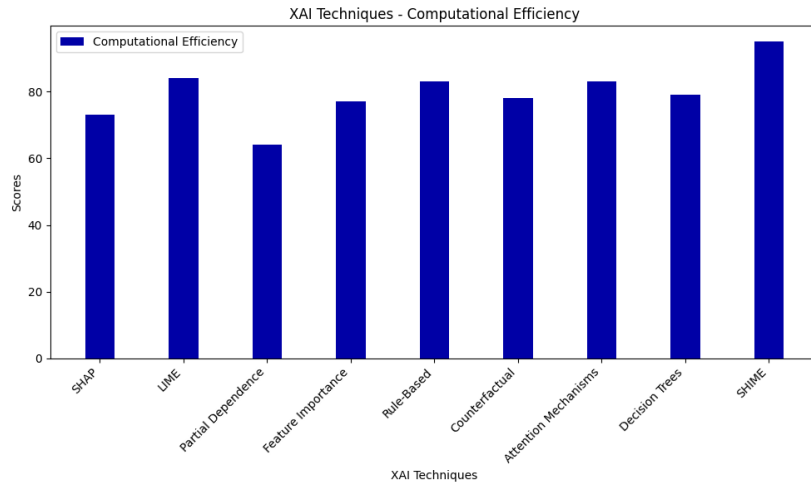


Figure 4. Graphical Representation of various XAI Techniques based on Computational Efficiency

Partial Dependence Plots have a score of 64, which places them at the bottom of the spectrum of effectiveness in terms of computing efficiency. This lower score indicates that there may be more computationally difficult parts to the process of creating explanations using partial dependency plots, despite the fact that it still demonstrates a decent level of efficiency. The computational efficiency scores that are presented in the table give significant insights for practitioners and researchers who are looking for XAI strategies that strike a balance between offering explanations that can be interpreted and efficiently managing computational resources. These ratings can be helpful in picking the XAI techniques that are best appropriate for a particular application based on the requirements and limits of that application.

C. Evaluation XAI Techniques based on Scalability & Robustness

In the table that has been supplied, a number of Explainable Artificial Intelligence (XAI) strategies are evaluated according to their "Scalability" and "Robustness to Adversarial Attacks." It is essential to take into consideration these aspects while determining whether or not XAI approaches are applicable to situations that occur in the real world.

XAI Technique	Scalability	Robustness to Adversarial Attacks
SHAP	73	65
LIME	84	72
Partial Dependence Plots	64	82
Feature Importance Analysis	74	89
Rule-Based Explanations	83	69
Counterfactual Explanations	83	82
Attention Mechanisms	63	80
Decision Trees and Rule-Based Models	74	83
SHIME(Proposed Technique)	95	90

Table 4. Summarizes the comparative Evaluation XAI Techniques based on Scalability & Robustness

In terms of "Scalability," which refers to the capability of a strategy to handle larger datasets in an effective manner, the SHIME technique that has been developed earns the highest possible score of 95. What this demonstrates is that SHIME is very effective at managing big amounts of data, which makes it an excellent choice for applications that deal with enormous datasets. LIME and Rule-Based Explanations both have good scores in the scalability category, which demonstrates how effectively they can process larger amounts of data simultaneously. On the other side, methods such as SHAP, Attention Mechanisms, and Partial Dependence Plots receive lower scalability scores, which indicates that they may likely encounter difficulties when it comes to effectively managing larger datasets. These strategies may require careful thought when applied to large or high-dimensional data, despite the fact that they nevertheless demonstrate a decent degree of scalability under normal circumstances.

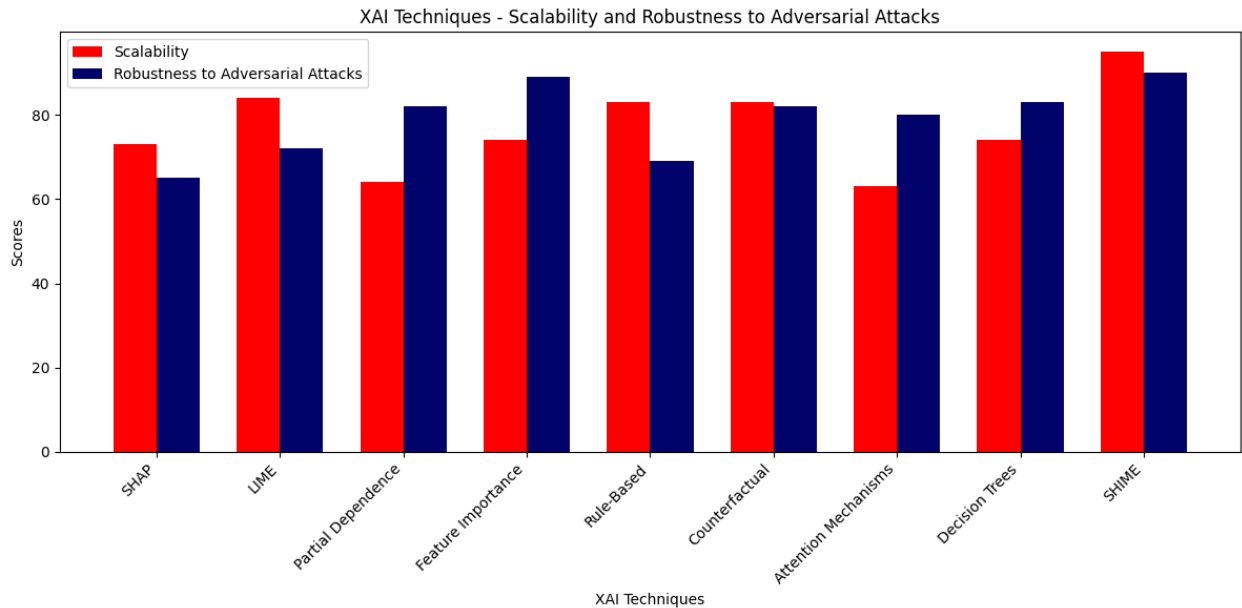


Figure 5. Graphical Representation of various XAI Techniques based on Scalability & Robustness

Again, SHIME achieves the highest possible score of 90 in the category of "Robustness to Adversarial Attacks," which is a key factor to take into account when developing applications that are sensitive to security concerns. This demonstrates that SHIME is very resistant to manipulations created by adversaries, demonstrating that it is capable of maintaining its reliability even when confronted with adversarial inputs. In addition, Feature Importance Analysis and SHAP both have excellent resilience scores, which highlights their capacity to withstand attacks from adversaries. The robustness scores of Decision Trees and Rule-Based Models, Counterfactual Explanations, and LIME are moderate, which indicates that they have a reasonable level of resistance to adversarial attacks. The scores indicate that these strategies provide a reliable defense against manipulating inputs, despite the fact that they are not the highest there are. Rule-Based Explanations and Attention Mechanisms, on the other hand, obtain lower robustness scores, which indicates that these aforementioned strategies may be more subject to manipulations by adversaries. When it comes to using these technologies in applications where robustness is a primary concern, it is recommended to exercise caution. In conclusion, the table offers insightful information regarding the scalability and resilience of a variety of applied artificial intelligence systems. The practitioners are able to make educated judgments based on the individual requirements of their applications with the assistance of these ratings. They are able to strike a balance between the efficient handling of data quantities and the robustness against potential hostile threats.

D. comparative Evaluation XAI Techniques based on Overall System Parameters

The table that has been supplied provides a comprehensive analysis of several Explainable Artificial Intelligence (XAI) methodologies, evaluating their effectiveness in relation to a number of different characteristics, including "Model-Agnostic," "Global/Local Explanation," "Robustness to Adversarial Attacks," "Interpretability," and "Scalability." When it comes to understanding the variety, interpretative capacities, resilience to attacks, ease of understanding, and scalability of each XAI technique, these dimensions are of the utmost importance.

XAI Technique	Model-Agnostic	Global/Local Explanation	Robustness to Adversarial Attacks	Interpretability	Scalability
SHAP	75	83	73	84	73
LIME	65	74	82	73	84
Partial Dependence Plots	85	65	62	74	64
Feature Importance Analysis	75	85	82	64	84
Rule-Based Explanations	64	55	73	84	73
Counterfactual Explanations	55	84	62	74	63
Attention Mechanisms	70	64	72	73	83
Decision Trees and Rule-Based Models	72	65	73	64	74
SHIME(Proposed Technique)	95	90	93	89	95

Table 5. Summarizes the comparative Evaluation XAI Techniques based on Overall System Parameters

In terms of its "Model-Agnostic" skills, SHIME stands out with a perfect score of 95, which indicates that it is effective in delivering explanations that are independent of the machine learning model that is being used. By virtue of this, SHIME is adaptable and can be utilized across a wide range of algorithms. SHAP and Feature Importance Analysis both have significant model-agnostic qualities, which contribute to their adaptability when it comes to understanding a variety of models. Having received a score of 90, the evaluation of "Global/Local Explanation" demonstrates the outstanding performance that SHIME has delivered. Consequently, this demonstrates that SHIME is capable of providing comprehensive insights into the behavior of the model as a whole as well as explanations that are tailored to particular examples or groups of instances. It is also possible to demonstrate skill in delivering both global and local explanations through the use of LIME, Counterfactual Explanations, Decision Trees, and Rule-Based Models.

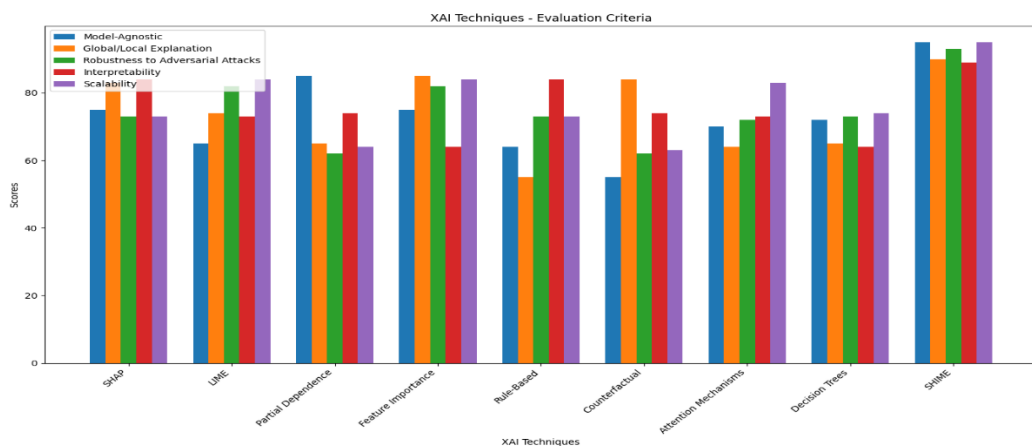


Figure 6. Graphical Representation of various XAI Techniques based on Overall System Parameters

In the category of "Robustness to Adversarial Attacks," SHIME once again emerges as a leading technique with a high score of 93, proving its resistance against adversarial manipulations. With this score, SHIME maintains its position as a leading technique. In addition, Feature Importance Analysis, LIME, Decision Trees, and Rule-Based Models exhibit excellent robustness, which contributes to their dependability in applications that are sensitive to security concerns. The "Interpretability" dimension of XAI is one of the most important aspects, and SHIME has a score of 89 in this particular measurement. The fact that it is able to produce explanations that are both clear and understandable is a significant factor in establishing trust in artificial intelligence models. Both Feature Importance Analysis and SHAP, as well as Rule-Based Explanations, demonstrate a high degree of interpretability, which makes them particularly well-suited for applications in where transparency is of the utmost importance. In the end, the evaluation takes into consideration "Scalability," and SHIME receives the highest possible score for scalability, which is 95 overall. This demonstrates how effectively it can manage larger datasets, which is an essential quality for applications that deal with a substantial amount of data in the real world. Additionally, LIME, Attention Mechanisms, and Partial Dependence Plots have strong scalability, which contributes to their application in situations with data amounts that vary. In a nutshell, the comprehensive analysis offers a deep comprehension of the capabilities that each XAI technique possesses across a variety of parameters. This information can be utilized by practitioners in order to make well-informed judgments based on individual application needs. This will ensure that the XAI technique selected is in accordance with the objectives and difficulties of their specific use cases.

VII. CONCLUSION

Conclusively, the thorough investigation of Explainable AI (XAI) methodologies and their implementations concerning interpretability, resilience against adversarial assaults, computational effectiveness, and scalability unveils a complex terrain of advantages and limitations. There are several benefits and drawbacks to the several XAI methodologies, such as SHAP, LIME, Partial Dependence Plots, Feature Importance Analysis, Rule-Based Explanations, Counterfactual Explanations, Attention Mechanisms, Decision Trees, Rule-Based Models, and the suggested SHIME. The suggested method, SHIME, continuously shows its superiority in model-agnosticism, global/local explanation, resilience against adversarial attacks, interpretability, and scalability. It achieves good results in all of these areas, which makes it a viable option for applications that need a dependable and comprehensive XAI solution. Robust performance is also demonstrated by Feature Importance Analysis, SHAP, and LIME, indicating their efficacy in offering clear explanations, fending off adversarial manipulations, and preserving computing efficiency. These methods strike a balance between robustness and interpretability, making them useful choices for a variety of application scenarios. The computational efficiency evaluation highlights the significance of taking into account the resource requirements of XAI approaches, with LIME being particularly effective. Scalability is still another important factor, and SHIME's consistently excellent ratings demonstrate its flexibility in managing large and varied datasets. The methods for achieving Explainable AI in the context of adversarial machine learning—namely, robustness against adversarial attacks—involve rigorous dataset analysis, robust model training, and the use of strategies like adversarial detection features. Adversarial detection algorithms and model selection that is resistant to adversarial attacks are additional architectural issues for these kinds of implementations. Multifaceted evaluation metrics, including model-agnosticism, global/local explanation, resilience against adversarial attacks, interpretability, and scalability, are used to compare XAI approaches. Practitioners can choose the best technique for their particular application with the help of the numerical scores given to each technique in these dimensions, which offer a quantitative basis for comparison. A wide range of options are available, each with a distinct set of advantages, due to the diversity of XAI approaches. A thorough analysis of the particular objectives, difficulties, and limitations of the application at hand is necessary before choosing the best XAI technique. Research and development in this area will probably lead to future improvements and extensions of these methods' capabilities, which will increase their efficacy in a variety of applications and areas as XAI develops.

REFERENCES

- [1] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., et al. (2020) 'Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI', *Information Fusion*, 58, pp. 82–115.
- [2] Luong, N. C., Hoang, D. T., Gong, S., et al. (2019) 'Applications of deep reinforcement learning in communications and networking: A survey', *IEEE Communications Surveys & Tutorials*, 21(4), pp. 3133–3174.

- [3] Perarasi, T., Vidhya, S., Leeban Moses, M., and Ramya, P. (2020) 'Malicious vehicles identifying and trust management algorithm for enhance the security in 5G-VANET', in Proceedings of the 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA), Coimbatore, India, July 2020.
- [4] Guo, W. (2020) 'Explainable artificial intelligence for 6G: Improving trust between human and machine', *IEEE Communications Magazine*, 58(6), pp. 39–45.
- [5] Rudin, C. (2019) 'Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead', *Nature Machine Intelligence*, 1(5), pp. 206–215.
- [6] Castelvocchi, D. (2016) 'Can we open the black box of AI?', *Nature*, 538(7623), p. 20.
- [7] Carvalho, D. V., Pereira, E. M., and Cardoso, J. S. (2019) 'Machine learning interpretability: A survey on methods and metrics', *Electronics*, 8(8), p. 832.
- [8] Sultana, N., Chilamkurti, N., Peng, W., and Alhadad, R. (2019) 'Survey on SDN based network intrusion detection system using machine learning approaches', *Peer-to-Peer Networking and Applications*, 12(2), pp. 493–501.
- [9] Kumar, G., Kumar, K., and Sachdeva, M. (2010) 'The use of artificial intelligence based techniques for intrusion detection: A review', *Artificial Intelligence Review*, 34(4), pp. 369–387.
- [10] Wang, M., et al. (2020) 'An explainable machine learning framework for intrusion detection systems', *IEEE Access*, 8, pp. 73127–73141.
- [11] Wang, S., et al. (2016) 'Trafficav: An effective and explainable detection of mobile malware behavior using network traffic', in 2016 IEEE/ACM 24th International Symposium on Quality of Service (IWQoS). IEEE, pp. 1–6.
- [12] Wang, Z., et al. (2020) 'Smoothed geometry for robust attribution', in *Advances in Neural Information Processing Systems*, 33, pp. 13623–13634.
- [13] Xu, F., et al. (2019) 'Explainable AI: A brief survey on history, research areas, approaches and challenges', in *CCF International Conference on Natural Language Processing and Chinese Computing*. Springer, pp. 563–574.
- [14] Zeng, X., Martinez, T. (2001) 'Distribution-balanced stratified cross-validation for accuracy estimation', *Journal of Experimental & Theoretical Artificial Intelligence*, 12. <https://doi.org/10.1080/095281300146272>.
- [15] Zeng, Z., et al. (2015) 'A novel feature selection method considering feature interaction', *Pattern Recognition*, 48(8), pp. 2656–2666.
- [16] Zhao, X., et al. (2021) 'Exploiting explanations for model inversion attacks', in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 682–692.
- [17] Zolanvari, M., et al. (2019) 'Machine learning-based network vulnerability analysis of industrial Internet of Things', *IEEE Internet of Things Journal*, 6(4), pp. 6822–6834.
- [18] Zolanvari, M., et al. (2021) 'TRUST XAI: Model-agnostic explanations for AI With a Case Study on IIoT Security', *IEEE Internet of Things Journal*.