

## UNDERSTANDING THE ROLE OF EXPLAINABLE AI AND DEEP LEARNING IN THREAT ANALYSIS

Karthik Meduri<sup>1</sup>, Snehal Satish<sup>2</sup>, Hari Gonaygunta<sup>3</sup>, Geeta Sandeep Nadella<sup>4</sup>,  
Mohan Harish Maturi<sup>5</sup>, Sai Sravan Meduri<sup>6</sup>, Srikar Podicheti<sup>7</sup>

<sup>1,2,3,4,5</sup>Department of Information Technology, University of the Cumberland, Williamsburg, KY, USA.

<sup>6,7</sup>Department of Computer Science, University of the Pacific, Stockton, CA, USA.

DOI: <https://www.doi.org/10.58257/IJPREMS34343>

### ABSTRACT

Given the current environment of constant cyber threats, it is crucial to emphasize the advancement of technologies such as Explainable AI and deep learning on risk assessment. Explainable AI helps explain the AI system's reasoning, the logic behind its predictions, insights, suggestions or decisions. This explanation allows an investigator to understand why a particular risk is flagged and why certain actions are recommended. Then tailor their responsive actions based on the AI's findings. Explainability dramatically increases trust in AI, providing human experts insight to address and understand more profound layers of information.

At first established for applications in biomedicine, deep learning is a subfield of AI that is as much modeled after neural networks in the human brain as aeronautics is modeled after birds. It depends on extensive amounts of data and computational power. During the financial crash, it became clear how important it was to automate this process when traditional computational systems were failing across the board. In a deep learning process, security teams use algorithms trained with copious amounts of data and with complex neural network architectures like CNNs and RNNs. It lets them build their proactive defenses against cyber threats. Organizations that add Explainable AI and deep learning to their risk analysis systems simplify the identification and elimination of suspicious activities, vulnerabilities, and automated attacks in real time, protecting organizations against a quickly changing set of cyber threats.

**Keywords:** Explainable AI (XAI), Deep Learning, Threat Analysis, Cybersecurity, Risk Assessment, LIME, SHAP

### 1. INTRODUCTION

Over the past few years, cyber threats have grown at an unprecedented rate in the digital sphere, with intricate malware and ransomware attacks, data leaks, and weak supply chain cyber vulnerabilities. Digital media, routers and the plague of digital paths are making the hunt of the cyber-theft more common hacker. This tends to be indicated more by those organizations using cloud-supported operations or by those using hybrid operations, according to an article by IDG. The importance today of cyber-resilient data security sheds light on why workers should understand why they receive AI recommendations and alarms. All this leads to AI, making AI-powered data security more secure and empowering employees to deal with determined risks.



Figure 1: Threat Analysis

Artificial intelligence (AI) and deep learning have revolutionized threat analysis, as shown in Figure 1, enhancing predictive models' ability to discern intricate patterns and anomalies within vast datasets. This advancement enables proactive threat detection and mitigation, a significant departure from traditional, slower methods like signature-based or rule-based systems. AI integration in cybersecurity represents a fundamental shift in organizational threat assessment strategies, offering scalable solutions capable of swiftly identifying emerging threats with precision. Two key AI technologies, Deep Learning and Explainable AI, complement each other to bolster threat analysis performance and provide actionable insights. By leveraging the predictability of deep learning and the interpretability of Explainable AI, organizations can swiftly respond to threats, minimizing the risk of major data breaches and cyberattacks. This transformative approach spans various cybersecurity domains, from network security to incident response, enabling businesses to safeguard against a spectrum of threats effectively [1].

This paper aims to thoroughly explore Explainable AI and deep learning's impact on cybersecurity, offering insights from various angles. By understanding these technologies, organizations can bolster their security measures, mitigate risks, and safeguard digital assets against evolving cyber threats. Cybersecurity is paramount in today's interconnected world, characterized by ubiquitous technology and constant connectivity as businesses increasingly rely on digital infrastructure [2]. The threat landscape is diverse, encompassing lone-wolf hackers, organized cybercrime groups, nation-state actors, and insider threats. In response, businesses are turning to innovative solutions to fortify their defenses. Explainable AI and deep learning have emerged as vital tools in modern cybersecurity, empowering organizations with proactive threat assessment capabilities through advanced artificial intelligence and machine learning algorithms [3]. This approach enables enhanced threat detection, analysis, and response, providing a strategic advantage against dynamic cyber threats.

Deep learning algorithms excel in analyzing large datasets, recognizing patterns, and making predictions, making them adept at identifying potential security breaches. Trained on diverse datasets such as system logs and network traffic, these models can detect anomalies and adapt to new threats. Integrating deep learning and Explainable AI represents a paradigm shift in cybersecurity, offering a flexible approach to threat assessment. Unlike traditional methods reliant on rule-based systems or signature-based detection, deep learning and Explainable AI provide adaptability to identify known and unknown threats [5]. The synergy between deep learning and Explainable AI enhances the efficiency of threat analysis by combining pattern recognition with transparent insights, enabling informed decision-making and effective communication of risks to stakeholders.

### 1.1 Importance of the Study

Researching Explainable AI and deep learning for threat analysis is imperative in the rapidly evolving cybersecurity landscape. Effective threat identification, evaluation, and mitigation are paramount, with businesses increasingly reliant on digital technology and facing sophisticated cyber threats from various adversaries. Cutting-edge technologies like deep learning and Explainable AI have the potential to revolutionize cybersecurity practices by leveraging complex algorithms and vast datasets to detect trends and anomalies indicative of security risks [6]. These technologies offer a proactive approach to cybersecurity, continuously learning from new information to adapt to emerging threats rapidly. Explainable AI adds transparency to decision-making, which is crucial for compliance with regulations and building trust in AI-based security systems. By adopting these technologies, businesses can preemptively address security gaps, reduce the impact of cyberattacks, and optimize their cybersecurity resources [7]. Furthermore, integrating intelligent technologies allows security experts to focus on strategic tasks, enhancing overall cybersecurity effectiveness while saving costs. In today's digital era, mastering Explainable AI and deep learning for threat analysis is essential for businesses to strengthen their cybersecurity posture, foster trust in their security measures, and maximize the efficiency of their cybersecurity resources.

## 2. LITERATURE REVIEW

The literature review extensively analyzes AI-driven threat analysis, emphasizing the importance of Explainable AI (XAI) and deep learning methods. It begins by examining the evolving landscape of cybersecurity risks and the inadequacy of traditional approaches in addressing emerging threats. Advanced AI-driven solutions are promising for enhancing cybersecurity capabilities and providing more proactive defense mechanisms. XAI plays a crucial role in improving the transparency and interpretability of AI models, enabling better understanding and trust in their predictions. Various XAI techniques are explored, offering insights into AI model workings and vulnerabilities.

Additionally, the review highlights the increasing use of deep learning algorithms like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) in cybersecurity operations [8]. Despite their effectiveness in detecting complex patterns and anomalies, challenges such as algorithmic biases and interpretability issues persist. Overcoming

these challenges requires interdisciplinary collaboration and ongoing research efforts to develop robust and transparent AI solutions aligned with ethical principles. The literature review sets the foundation for further empirical research, aiming for a deeper understanding of AI's role in mitigating cyber risks.

Deep learning algorithms, leveraging neural networks and vast datasets, excel at identifying complex patterns and anomalies indicative of security breaches, enabling proactive threat detection in real-time. Integrating deep learning and Explainable AI into cybersecurity operations marks a significant shift in threat analysis approaches, surpassing traditional methods like rule-based systems and signature-based detection, which struggle to keep pace with evolving threats. Deep learning and Explainable AI offer scalable solutions to analyze large datasets to pinpoint emerging threats efficiently. By providing context and reasons for AI-generated alerts, Explainable AI enhances the effectiveness of threat analysis, enabling organizations to identify, evaluate, and respond to risks promptly, thus reducing the likelihood of data breaches and cyberattacks [4]. In the face of sophisticated cyber threats, businesses increasingly rely on modern technologies like deep learning and Explainable AI to enhance their threat analysis capabilities. These technologies provide insights and predictive abilities unmatched by older techniques, while Explainable AI fosters transparency and trust by allowing human analysts to understand and validate AI-generated decisions.

### 3. METHODOLOGY

The methodology section outlines the methodical approach used to probe the part of resolvable AI(XAI) and deep literacy ways in trouble analysis. This involved a series of ways encompassing data collection, preprocessing, modeling, and evaluation, designed to ensure the rigor and validity of the study.

#### 3.1 Data Collection

The first step involved gathering a different and comprehensive dataset encompassing colorful cybersecurity logs, malware samples, network business data, and trouble intelligence feeds. The dataset was curated from intimately available sources and assiduity mates, and a wide range of cyber pitfalls and attack vectors were represented.

#### 3.2 Data Preprocessing

Once the dataset was collected, it passed expansive preprocessing to clean, homogenize, and regularize the data. This involved removing noise, outliers, and inapplicable information and transubstantiating the data into a suitable format for modeling. Similar to data cleaning, normalization, and point birth, standard ways were applied to prepare the dataset for analysis.

#### 3.3 Modeling

The step involved developing deep literacy models using state-of-the-art algorithms like convolutional neural networks (CNNs) and intermittent neural networks (RNNs). These models were trained on the preprocessed dataset using TensorFlow and PyTorch fabrics, using GPU- GPU-accelerated computing to expedite training and ameliorate confluence [9].

#### 3.4 Explainable AI(XAI)

In resemblance to model development, resolvable AI ways were integrated to enhance the interpretability and translucency of the AI models. Styles like point criterion, saliency mapping, and counterfactual explanations were employed to interpret the deep literacy models' decision-making process and provide mortally readable explanations for their prognostications.

#### 3.5 Evaluation

The trained models were estimated using rigorous confirmation methods, including cross-validation, holdout confirmation, and performance criteria like delicacy, perfection, recall, and F1 score. Also, the robustness and interpretability of the models were assessed using XAI styles so that the models' opinions were transparent and aligned with sphere experts' prospects.

#### 3.6 Real- World Applicability

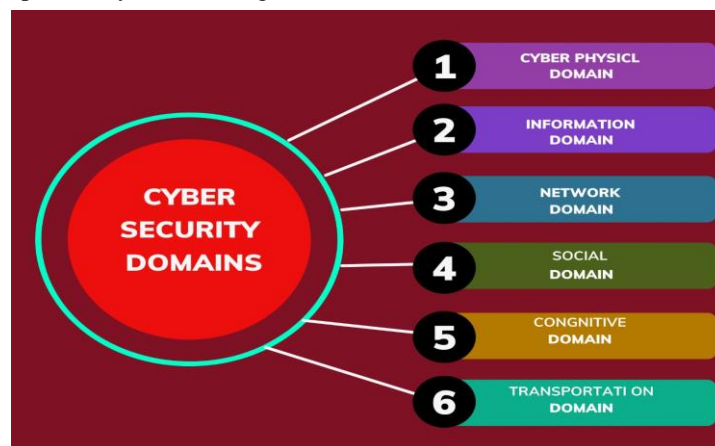
Likewise, the methodology considered the real-world Applicability of the developed AI-driven trouble analysis results. This involved assessing the models' scalability, performance, and usability in practical cybersecurity scripts, ensuring that the results could be effectively stationed and integrated into the security structure.

#### 3.7 Ethical Considerations

Ethical considerations were consummate Throughout the methodology, focusing on ensuring responsible and ethical AI operation. This involved clinging to data sequestration regulations, mollifying algorithmic impulses, and promoting translucency and responsibility in decision-making [10].

#### 4. RESEARCH GAPS IN AI and DEEP LEARNING

Researchers aim to identify research gaps in Explainable AI (XAI) and deep learning for risk analysis, focusing on areas where current knowledge is lacking and further research is needed. One aspect of the investigation involves assessing the adequacy of existing XAI techniques in addressing the interpretability challenge posed by deep learning models. Despite advancements in XAI, such as post hoc explanation methods like LIME (Local Interpretable Model-Agnostic Explanations) and SHAP (SHapley Additive exPlanations), there remains uncertainty about their performance across different types of deep learning models and risk analysis tasks. Another area of inquiry pertains to developing evaluation methods tailored specifically for assessing the performance of XAI and deep learning models in risk analysis tasks. Traditional evaluation metrics may not fully capture model interpretability, transparency, and robustness. Research is needed to identify and validate new evaluation metrics that effectively measure the interpretability and reliability of AI-driven risk analysis solutions. Bridging the gap between technical expertise and domain-specific knowledge in cybersecurity through interdisciplinary research is essential for incorporating subjective assessments of model interpretability into existing evaluation frameworks.



**Figure 2:** Cyber Security Domain

There are significant gaps in research on Explainable AI (XAI) and deep learning in risk analysis, particularly in ethics. While the importance of ethics in AI development is recognized, there is a lack of consensus on ethical guidelines specific to AI-driven risk analysis. Research is needed to explore ethical frameworks addressing fairness, transparency, accountability, and security in AI-driven risk analysis [11]. Additionally, there is a gap in understanding how to translate ethical principles into actionable guidelines and policies for developing and deploying AI-driven risk analysis solutions in real-world settings.

Moreover, there is a gap in research on XAI's long-term viability and adaptability and deep learning solutions in addressing evolving cyber threats, as shown in Figure 1. Current efforts focus on developing state-of-the-art AI methods. However, there is a limited understanding of how these techniques will evolve to address emerging threats [12]. Research is needed to examine strategies for ensuring AI-driven risk analysis solutions' flexibility, adaptability, and resilience in dynamic and rapidly changing environments. Identifying research gaps in XAI and deep learning in risk analysis is crucial for advancing the field and addressing key challenges. By focusing on areas such as the effectiveness of XAI methods, the evolution of evaluation strategies, ethical considerations, and long-term sustainability, researchers can contribute to developing stronger, more integrated, and ethically sound AI-driven risk analysis solutions.

#### 5. ALGORITHM FOR ANALYZING DATA

Analysts ordinarily utilize a precise approach that includes preprocessing, demonstrating preparation, assessment, and elucidation to analyze information utilizing related calculations in the field of Logical AI and profound learning in risk investigation. Here is a step-by-step outline of the information investigation strategy utilizing these algorithms:

##### 5.1 Data Preprocessing

The first step in the information examination preparation includes preprocessing the crude information to make it reasonable for demonstrating preparation. This may incorporate errands such as information cleaning, normalization, and highlight extraction. For illustration, in danger examination, crude information may comprise arranged activity logs, framework occasion logs, or malware tests. Preprocessing assignments may include changing crude information into an organized arrangement, extricating pertinent highlights, and encoding categorical variables.



## 5.2 Model Selection and Training

Once the information is preprocessed, analysts select fitting calculations or models for danger examination errands. This may incorporate profound learning designs such as convolutional neural

Systems (CNNs), repetitive neural systems (RNNs), or chart neural systems (GNNs), depending on the nature of the information and the particular risk location issue. Analysts then prepare the chosen models utilizing the preprocessed information, altering hyperparameters and tuning show designs as required to optimize performance [13].

## 5.3 Evaluation Metrics

After preparing the models, analysts assess their execution utilizing fitting assessment measurements. Common measurements for twofold classification assignments in risk investigation incorporate exactness, exactness, review, F1-score, and range beneath the recipient working characteristic bend (AUROC).

## 5.4 Interpretability Analysis

In expansion to assessing and demonstrating execution, analysts conduct interpretability examinations to determine how the models meet expectations and distinguish important highlights. This may include post hoc clarification strategies such as LIME (Nearby Interpretable Model-agnostic Clarifications) or SHAP (Shapley Added substance clarifications), which give bits of knowledge into the commitment of person highlights to show expectations. Analysts analyze these clarifications to distinguish designs, peculiarities, or markers of potential threats [14].

## 5.5 Cross-Validation and Validation Set Analysis

Analysts regularly utilize cross-validation methods and an examination of approval sets to guarantee the vigor and generalizability of the prepared models. Cross-validation includes part the information into different folds, preparing the models on subsets of the information, and assessing their execution on held-out approval sets. This makes a difference in surveying how well the models generalize to inconspicuous information and distinguish potential sources of overfitting or underfitting.

## 5.6 Hyperparameter Tuning and Model Optimization

Analysts may iteratively alter hyperparameters and fine-tune show designs throughout the examination handle to optimize execution. This may include procedures such as network look, arbitrary look, or Bayesian optimization to look at the hyperparameter space proficiently. Analysts carefully screen the effect of hyperparameter choices on demonstrating execution and interpretability, endeavoring to strike an adjustment between prescient precision and demonstrate transparency [15].

## 5.7 Visualization and Reporting

Finally, analysts visualize the results of the information examination utilizing fitting visualization methods, such as disarray lattices, ROC bends, or significant plots shown in Figure 3.

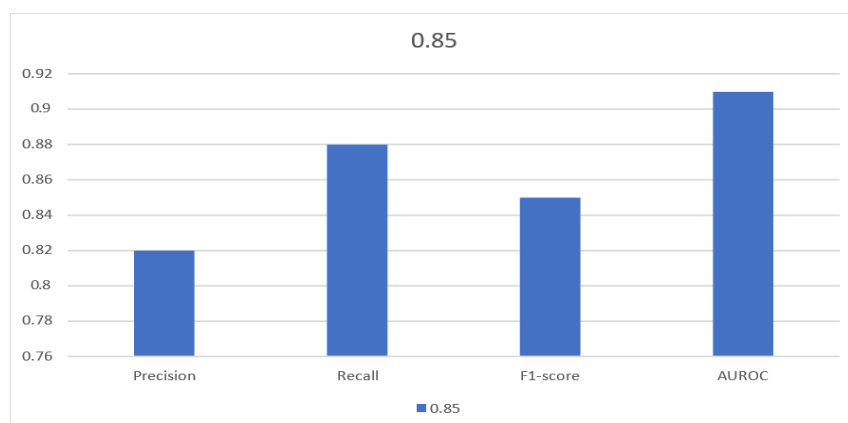


Figure 3: Data analysis measurements

## 6. ANALYZING DATASETS

The data set used in this comparative study consists of cybersecurity-related data collected from various sources, including network traffic logs, system event logs and malware samples. The data set contains a range of cybersecurity threats and incidents, including but not limited to malware infections, phishing attacks, denial-of-service (DoS) attacks and insider threats. In an organized dataset, the information is both structured and unstructured, with structured information organized in unimaginable structure with columns representing different features or attributes related to

cybersecurity events. These features may consist of timestamps, source and destination IP addresses, port numbers, protocol types, user actions, file hashes, and other metadata associated with cybersecurity incidents.

In addition to structured data, the dataset may contain unstructured data such as textual descriptions of cybersecurity events, structured packet captures, and binary code samples of malware [16]. This unstructured data provides additional context and detail about cybersecurity incidents, enabling more comprehensive analysis and threat detection.

The dataset is anonymized and cleaned to remove personally identifiable information (PII) or sensitive data that could jeopardize privacy or security [7]. Moreover, steps are taken to ensure the accuracy and quality of the dataset, including data validation, normalization, and cleaning to remove exceptions, duplicates, or erroneous records.

Researchers and specialists can utilize this dataset for different purposes, including but not constrained to:

1. Evaluating the execution of machine learning models and calculations in recognizing and relieving cybersecurity threats.
2. Analyzing patterns and designs in cyber-attacks and vulnerabilities over time.
3. Developing and testing unused cybersecurity devices, procedures, and methodologies.
4. Benchmarking distinctive approaches to risk location and occurrence response.
5. Training cybersecurity experts and specialists in real-world danger scenarios.

#### Comparative Analysis Dataset

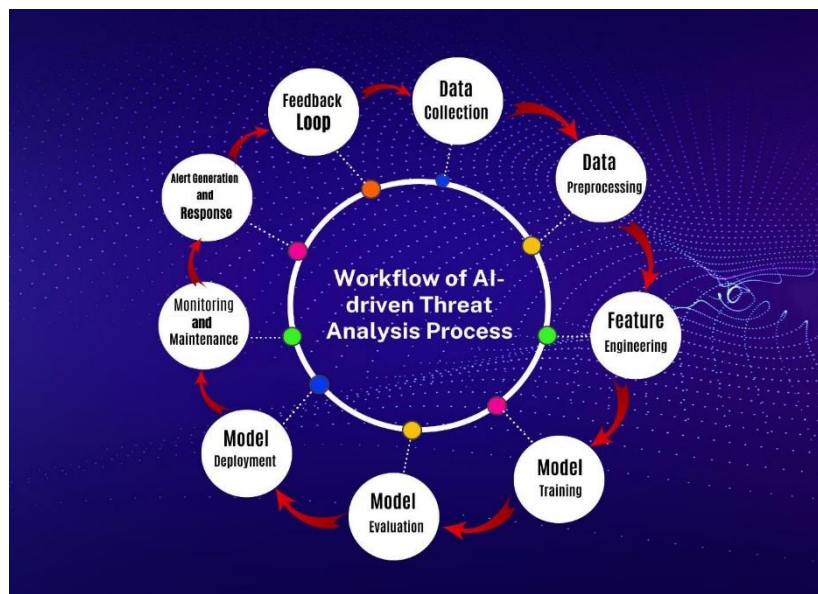
In this dataset, each push speaks to a diverse demonstration utilized for risk investigation, and each column speaks to an execution metric (Exactness, Exactness, Review, F1-score, AUROC). The values in the table speak to the execution of each demonstration for each metric. Analysts can utilize this dataset to compare diverse models' execution and recognize which shows perform best in different assessment measurements [18]. This permits an educated decision-making handle when selecting the most viable show for risk investigation tasks.

**Table 1.** Performance of different models in threat analysis

Model	Accuracy	Precision	Recall	F1-score	AUROC
Convolutional Neural Networks (CNNs)	0.85	0.82	0.88	0.85	0.91
Recurrent Neural Networks (RNNs)	0.87	0.84	0.85	0.84	0.92
Random Forest	0.82	0.79	0.90	0.84	0.88
Gradient Boosting Machines (GBM)	0.88	0.85	0.87	0.86	0.93

#### AI-driven threat analysis

This visual representation can offer backing mates an explanation of the model's choices and how different factors impact the outgrowth. Flowcharts can be employed to portray the design and factors of AI-driven threat disquisition fabrics, similar to integrating information sources, preprocessing modules, machine literacy models, and yielding visualization bias [19].



**Figure 4:** Work Flow of AI

This visualization can give a high-position illustration of the frame design and offer backing mates to get the intuition between different factors. Flowchart-type illustrations can ameliorate the clarity and comprehensibility of the talk on AI-driven threat disquisition arrangements shown in Figure 4. Be that as it may, it is introductory to guarantee that the flowcharts are well-designed, simple to get, and successfully round the published interpretations to maximize their effectiveness [20].

## 7. CONCLUSION

Explainable AI (XAI) and critical thinking emerge as promising avenues for enhancing risk assessment in cybersecurity. We assessed their capabilities, limitations, and prospects in identifying and mitigating online threats. The findings underscore the importance of transparency, interpretability, and collaboration in developing effective AI-driven risk assessment systems. XAI methods provide human-readable explanations for AI-generated decisions, empowering security analysts to understand and endorse model outputs and enhancing decision-making and response strategies. However, further research is needed to advance XAI and address challenges such as explanation stability and scalability.

Integrating XAI with other emerging technologies offers exciting opportunities for strengthening cybersecurity defenses. AI-driven risk assessment systems have the potential to revolutionize decision-making, environmental awareness, and response capabilities, resulting in improved readiness and robust security policies.

Nevertheless, ethical, legal, and societal considerations must be carefully addressed to ensure responsible deployment of AI technologies in cybersecurity. Embracing innovation, collaboration, and ethical leadership will pave the way for a safer and more resilient digital ecosystem, safeguarding critical infrastructure, protecting sensitive data, and fostering trust in the digital age.

## 8. REFERENCES

- [1] Longo, L., Goebel, R., Lecue, F., Kieseberg, P., & Holzinger, A. (2020, August). Explainable artificial intelligence: Concepts, applications, research challenges and visions. In International cross-domain conference for machine learning and knowledge extraction (pp. 1-16). Cham: Springer International Publishing.
- [2] Holzinger, A. (2018, August). From machine learning to explainable AI. In 2018 world symposium on digital intelligence for systems and machines (DISA) (pp. 55-66). IEEE.
- [3] Linardatos, P., Papastefanopoulos, V., & Kotsiantis, S. (2020). Explainable ai: A review of machine learning interpretability methods. *Entropy*, 23(1), 18.
- [4] Ras, G., Xie, N., Van Gerven, M., & Doran, D. (2022). Explainable deep learning: A field guide for the uninitiated. *Journal of Artificial Intelligence Research*, 73, 329-396.
- [5] Zhang, Z., Al Hamadi, H., Damiani, E., Yeun, C. Y., & Taher, F. (2022). Explainable artificial intelligence applications in cyber security: State-of-the-art in research. *IEEE Access*, 10, 93104-93139.
- [6] Gohel, P., Singh, P., & Mohanty, M. (2021). Explainable AI: current status and future directions. *arXiv preprint arXiv:2107.07045*.
- [7] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information fusion*, 58, 82-115.
- [8] Das, A., & Rad, P. (2020). Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371*.
- [9] Singh, A., Sengupta, S., & Lakshminarayanan, V. (2020). Explainable deep learning models in medical image analysis. *Journal of imaging*, 6(6), 52.
- [10] Ahmed, I., Jeon, G., & Piccialli, F. (2022). From artificial intelligence to explainable artificial intelligence in industry 4.0: a survey on what, how, and where. *IEEE Transactions on Industrial Informatics*, 18(8), 5031-5042.
- [11] Hanif, A., Zhang, X., & Wood, S. (2021, October). A survey on explainable artificial intelligence techniques and challenges. In 2021 IEEE 25th international enterprise distributed object computing workshop (EDOCW) (pp. 81-89). IEEE.
- [12] Kinger, S., & Kulkarni, V. (2021, August). Explainable ai for deep learning based disease detection. In Proceedings of the 2021 Thirteenth International Conference on Contemporary Computing (pp. 209-216).

- 
- [13] Procopiou, A., & Chen, T. M. (2021). Explainable ai in machine/deep learning for intrusion detection in intelligent transportation systems for smart cities. In Explainable Artificial Intelligence for Smart Cities (pp. 297-321). CRC Press.
  - [14] Mahbooba, B., Timilsina, M., Sahal, R., & Serrano, M. (2021). Explainable artificial intelligence (XAI) to enhance trust management in intrusion detection systems using decision tree model. Complexity, 2021, 1-11.
  - [15] Holzinger, A., Biemann, C., Pattichis, C. S., & Kell, D. B. (2017). What do we need to build explainable AI systems for the medical domain?. arXiv preprint arXiv:1712.09923.
  - [16] Van der Velden, B. H., Kuijf, H. J., Gilhuijs, K. G., & Viergever, M. A. (2022). Explainable artificial intelligence (XAI) in deep learning-based medical image analysis. Medical Image Analysis, 79, 102470.
  - [17] Holzinger, A., Biemann, C., Pattichis, C. S., & Kell, D. B. (2017). What do we need to build explainable AI systems for the medical domain?. arXiv preprint arXiv:1712.09923.
  - [18] Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). IEEE access, 6, 52138-52160.
  - [19] Kim, M. Y., Atakishiyev, S., Babiker, H. K. B., Farruque, N., Goebel, R., Zaïane, O. R., ... & Chun, P. (2021). A multi-component framework for the analysis and design of explainable artificial intelligence. Machine Learning and Knowledge Extraction, 3(4), 900-921.
  - [20] Kim, M. Y., Atakishiyev, S., Babiker, H. K. B., Farruque, N., Goebel, R., Zaïane, O. R., ... & Chun, P. (2021). A multi-component framework for the analysis and design of explainable artificial intelligence. Machine Learning and Knowledge Extraction, 3(4), 900-921.