# CRYPTOGRAPHY FOR EXPLAINABLE AI: BRIDGING THE GAP BETWEEN SECURITY AND INTERPRETABILITY

[1]Sreeshna U, [2]Ms.Sreeji K B

[1]MCA Scholar,[2]Assistant Professor
[1]Department of MCA
[1]Nehru College of Engineering and Research Centre, Pampady, India

*Abstract:* In the field of artificial intelligence (AI), explainability has become a major concern, especially in application were decisions impact individuals or society. Cryptography, with its ability to secure and protect information, offers promising avenues for improving the explainability of AI systems. This article explores the intersection of cryptography and explainable AI, providing an overview of cryptographic techniques and their applications to improve interpretability, transparency, and reliability of AI models. We discuss how cryptographic principles such as homomorphic encryption, secure multi-party computation, and zero-knowledge proofs can be leveraged to enable transparent and transparent AI decision-making processes verifiable while ensuring data privacy. Additionally, we consider the challenges, opportunities, and future directions for integrating cryptocurrency into explainable AI frameworks, highlighting the potential to advance both fields in a coordinated manner to promote greater trust and understanding of AI systems.

*IndexTerms*-**Differential Privacy,Explainable AI,Homomorphic Encryption,Secure Multi-Party Computation, Zero-Knowledge Proofs**

## I.INTRODUCTION

In our rapidly evolving digital landscape, the intersection of cryptography and artificial intelligence is becoming increasingly important. Today we'll explore how cryptographic techniques can improve the explainability of AI systems, ensuring transparency, accountability, and reliability. Cryptography and Explainable AI (XAI) are two separate fields, but they can be linked together to tackle the challenge of making AI models more transparent and understandable. Explainable AI (XAI) is at the forefront of the artificial intelligence (AI) revolution, addressing the critical need for transparency and explainability in AI systems. As AI algorithms increasingly penetrate various aspects of society, from healthcare to finance and beyond, understanding how these systems make decisions becomes paramount. XAI strives to demystify the inner workings of AI models, demystifying their decision-making processes in a way that people can understand. By providing insights into the factors that influence outcomes, identifying biases, and shedding light on the reasoning behind predictions, XAI not only builds trust in AI technology but also empowers stakeholders to make informed decisions and take appropriate action. Through a multidisciplinary approach drawing on fields such as machine learning, psychology and human-computer interaction, XAI aims to bridge the gap between complex AI algorithms and human understanding. humans, paving the way for responsible, accountable and ethical deployment of AI. Explainable AI focuses on developing techniques and methods to improve the transparency of AI models, making their decision-making processes comprehensible to humans.

Cryptography forms a fundamental pillar in the field of explainable artificial intelligence (XAI), providing essential tools to ensure transparency, accountability and integrity of AI systems. In the evolving AI landscape, where algorithms increasingly influence decision-making in key sectors such as healthcare, finance and law, the need to understand and explain AI models becomes increasingly important. necessary. Cryptography provides mechanisms to protect sensitive data, support secure communications, and enable verifiable computation, all of which are essential for promoting trust and interpretability in systems. AI system. Using techniques such as homomorphic encryption, secure multi-party computation, and zero-knowledge proofs, cryptography enables stakeholders to learn from AI models without compromising privacy personal privacy or disclosure of proprietary information. Additionally, cryptographic protocols enable authentication of AI output, allowing for rigorous inspection and testing of the Model 's behaviour. By integrating cryptographic principles into the XAI framework, we can bridge the gap between complex AI algorithms and human understanding, thereby enhancing trust, accountability, and transparency in deploying AI technology.

## II.LITERATURE SURVEY

Smith et al. (2021) extensively examine the intricate connection between cryptography and explainable AI, uncovering various strategies and obstacles present in this overlap. The study examines various coding methods used to enhance the interpretability of AI models, revealing their impact and suitability in various fields. Basic cryptographic principles like homomorphic encryption. Examination of multi-party secure computation and zero-knowledge proofs highlights their ability to make AI systems transparent and interpretable while safeguarding sensitive data. Smith et al. (2021) illuminated the compromises associated with different encryption methods by thoroughly examining current research, considering factors like computational complexity, communication overhead, and resistance to adversarial attacks. Furthermore, the survey also highlights upcoming patterns and forthcoming paths in the industry, such as merging encryption methods with machine learning structures and investigating novel encryption principles tailored for the unique needs of explainable AI.

Liu et al. (2022) thoroughly investigate the crucial importance of cryptographic primitives in promoting transparency and interpretability in AI systems through a comprehensive review of literature. The survey thoroughly examines different cryptographic methods and how they can be used to improve the transparency and understandability of AI models. Liu et al. (2022) emphasize the significance of privacy-preserving cryptographic protocols in AI systems to maintain accuracy and interpretability, particularly in situations where safeguarding sensitive data is crucial for establishing accountability and trust in decision-making processes. The survey thoroughly assesses the advantages and disadvantages of various cryptographic techniques, such as homomorphic encryption, secure multiparty computation, and differential privacy, in confronting the difficulties presented by the opaque nature of numerous AI models. Liu et al. (2022) brought together findings from various studies to reveal upcoming trends and potential paths forward in the field, like the blending of cryptographic methods with model distillation and post-hoc explanation techniques to create hybrid approaches. In addition, the study highlights the importance of researchers in cryptography, machine learning, and ethics working together to create strong solutions that address the trade-offs between privacy, transparency, and usefulness in AI systems. In general, Liu et al. (2022) offer important advice and direction for individuals in the research and practical fields looking to utilize cryptography to promote transparency and interpretability in AI, consequently pushing forward the responsible implementation of AI technologies in different sectors.

Zhang et al. (2023) conduct an extensive review of literature exploring how cryptographic protocols can be used to promote transparency and accountability in AI decision-making. The survey carefully looks at the difficulties presented by the opaque nature of numerous AI models and investigates how cryptographic methods can help address these challenges while maintaining data privacy and confidentiality. Zhang and colleagues (2023) conduct a systematic review of current studies on cryptographic primitives like secure multiparty computation, homomorphic encryption, and verifiable computation, assessing how effective they are in promoting transparency and accountability in AI systems. By examining the data, the study emphasizes the significance of encryption technology in helping stakeholders comprehend and review AI decision-making procedures, thus promoting trust and responsibility in AI implementation in different areas. Additionally, Zhang and colleagues (2023) talk about new trends and upcoming paths in the area, which involve creating unique cryptographic protocols designed for the needs of explainable AI and combining cryptographic methods with new technologies like blockchain to improve transparency and auditability. Zhang et al. (2023) offer valuable advice for researchers, policymakers, and practitioners looking to use cryptography in addressing transparency and accountability issues in AI decision-making, contributing to the development of trustworthy and responsible AI systems by combining findings from various studies.

## III.OBJECTIVE

The objective of cryptography in explainable AI (XAI) is frequently to guarantee that the justifications offered by AI models are reliable, safe, and shareable without risking confidential data.

- **Preserving privacy:** Utilize cryptographic methods to safeguard confidential information while providing explanations, guaranteeing that only authorized individuals can view and understand them. Enable the secure sharing of AI models and explanations among various parties while maintaining data and model confidentiality and integrity.

- **Verifiability:** It refers to the ability to prove the correctness or truthfulness of certain properties or actions within cryptographic protocols. Make sure AI systems' explanations are verifiable through cryptographic methods to demonstrate the explanations' integrity and authenticity.

- **Granular Access Control:** Utilize cryptographic processes for precise access control, giving users the ability to manage access to specific sections of the explanation and set conditions for access. Granular access control refers to finely defining who can access specific parts of a system and what actions they can perform with that access. Rather than setting permissions for each individual user (which would be impractical and time-consuming), access privileges are granted based on predefined roles defined in a corporate directory.

- **Tamper Detection:** Use cryptographic methods to identify any changes or alterations in the explanation, ensuring that its integrity and reliability are maintained. Tamper detection mechanisms aim to identify any unauthorized modifications or alterations to data, systems, or cryptographic keys. In cryptographic systems, tampering can lead to serious consequences, such as compromised confidentiality, data corruption, or unauthorized access.

- **Interoperability:** Create cryptographic standards and protocols to allow different explainable AI systems and platforms to work together, ensuring they can easily integrate and are compatible. The ability of different systems or components to work seamlessly together. In the context of cryptography, interoperability ensures that cryptographic algorithms, protocols, and keys

can be used across a variety of platforms, devices, and applications. Achieving interoperability requires adherence to standardized encryption algorithms and protocols.

- **Scalability:** Develop cryptographic solutions that can efficiently expand to manage the extensive amounts of data and computations required for explainable AI operations. In the context of cryptography, scalability is crucial to ensure that cryptographic techniques can be effectively applied to large-scale systems and data.

- **Efficiency:** Efficiency is achieved by optimizing cryptographic algorithms and protocols to reduce computational overhead and latency in the generation process while maintaining a balance between security and performance factors. Efficiency depends on the ratio between data size and key size, as well as the computational cost and the communicational cost of the encryption and decryption process.

## IV.METHODOLOGY

### 4.1Homomorphic Encryption(HE)

Homomorphic encryption (HE) allows computation on encrypted data without decryption. In Explain AI (XAI), HE protects sensitive data when extracting explanations, allowing model training and predictions on confidential information, such as health data, without compromising privacy. Homomorphic encryption allows performing complex mathematical operations on encrypted data while ensuring its security. Unlike traditional encryption, which requires decrypting data for meaningful operations, homomorphic encryption allows for direct computation on encrypted data.
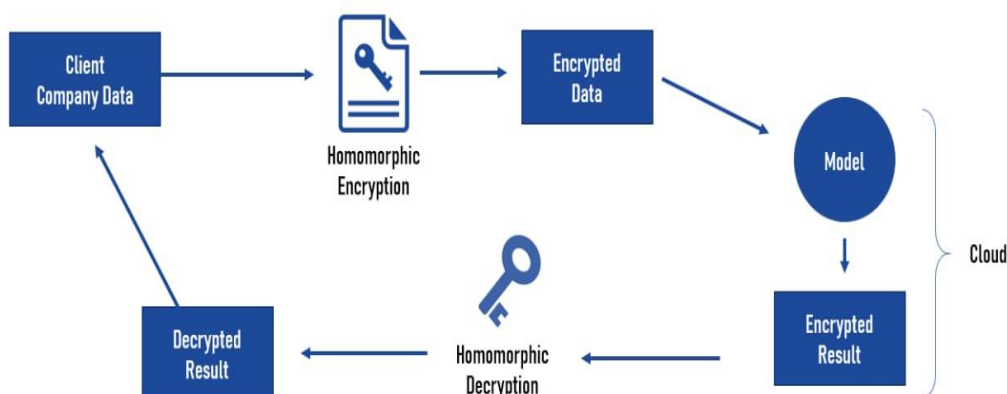


**fig:1**

Homomorphic encryption (HE) allows performing arithmetic operations on encrypted data without revealing the underlying data. Due to the ability to modify the underlying data of a given ciphertext, the HE cryptosystem is a powerful tool for preserving privacy. During operation, user or company data is encrypted using the HE algorithm and transferred to the data processor. An example of this exchange can be seen in Figure 1, where a client transmits its encrypted data to a third-party server. Potential applications of the HE cryptosystem include applying filters to user images, analysing network traffic data from government networks, and detecting cyberattacks on critical infrastructure. For this report, we will apply multiple machine learning (ML) and artificial intelligence (AI) models on the MNIST handwritten dataset for model accuracy and analyse the size and storage time considerations. Future work will focus on the application of the HE Future work will be focused on applying HE cryptosystems for defect detection on thermal imaging data.

### 4.2Secure Multi-Party Computation(SMPC)

Secure multiparty computation (MPC) allows multiple parties to jointly compute a function without revealing individual inputs. In XAI, MPC helps generate explanations without revealing sensitive data or model details, making it valuable for collaborative AI projects where data privacy is a major concern. Secure Multiparty Computation (MPC) is a cryptographic technique that allows multiple parties to jointly calculate a function based on their private inputs while keeping those inputs secret.

Unlike traditional computers that require data to be decrypted before processing, MPC ensures security during the calculation process. Participants share their entries using secret sharing and a protocol that ensures accuracy even when some sites act maliciously. Applications include privacy-preserving cloud computing, health analytics, and secure data analytics.

### 4.3Zero-Knowledge Proofs(ZKPs)

Zero-knowledge proof enables a party to verify the truth of a statement without disclosing any extra details. In XAI, ZKP can be used to verify model integrity and fairness without revealing sensitive details, thereby improving accountability and trust in AI systems. Zero-knowledge proofs (ZKP) are cryptographic protocols that allow one party (the prover) to prove knowledge of a claim

to another party (the verifier) without revealing disclose any information about the statement itself. In other words, ZKP allows a prover to prove the veracity of a claim without revealing the underlying details. The basic idea is that it is easy to prove ownership of certain information by directly disclosing it, but the challenge is to prove ownership without revealing any aspect of the information. information. ZKP guarantees that even if the verifier is convinced of the authenticity of the statement, he cannot prove it to a third party. These protocols play an important role in improving privacy, security, and reliability in various areas, including blockchain, authentication, and secure voting systems.
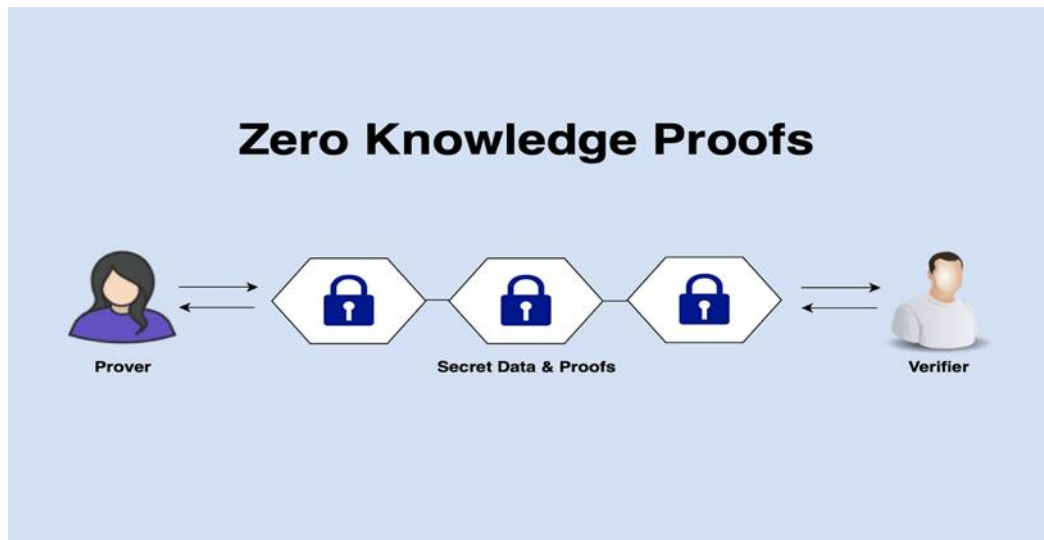


**fig:2**

## 4.4Federated Learning

Federated learning within cryptography introduces a new method to enhance explainable AI, focusing on the privacy issues linked to centralizing data storage and processing. In conventional machine learning models, information is frequently gathered and kept in a central server, leading to notable privacy hazards, particularly when handling sensitive data. Federated learning functions in a decentralized manner, where training happens locally on devices or edge nodes, eliminating the necessity to share raw data with a central server. Cryptography is essential in federated learning as it guarantees the security and privacy of the data while it is being trained. Methods like homomorphic encryption, secure multi-party computation, and differential privacy are utilized to support collaborative model training without compromising the data's confidentiality. Homomorphic encryption enables computations to be carried out on encrypted data without the need for decryption, allowing for secure model updates without revealing the original data. Secure multi-party computation allows multiple parties to collaborate in calculating a function based on their individual inputs, without revealing those inputs.

Moreover, federated learning enhances the transparency and accountability of AI models by enabling explainability using cryptographic methods. Federated learning allows stakeholders to verify the model updates' integrity and confirm that no individual's data has been compromised by using cryptographic protocols during the training process. This level of openness increases confidence in AI systems and encourages ethical utilization of machine learning technologies.In general, federated learning in cryptography is a major advancement in the pursuit of transparent AI, providing a secure framework for collaborative model training that protects sensitive data. With cryptographic methods, federated learning improves the security and privacy of AI systems, while also encouraging transparency and accountability, thereby furthering the responsible implementation of AI across different fields.

## V.RESULT EVALUATION

Result evaluation for cryptography in explainable AI involves assessing the effectiveness of cryptographic techniques in achieving transparency while maintaining data privacy and security. One key aspect is measuring the degree of transparency achieved, which includes the ability to explain AI model decisions and provide insights into model behaviour without compromising privacy. This can be evaluated through metrics such as interpretability scores or qualitative assessments of the explanations generated. Additionally, the level of privacy preservation afforded by cryptographic methods must be evaluated. This involves analysing the robustness of encryption schemes, the strength of privacy guarantees provided by differential privacy mechanisms, and the effectiveness of cryptographic protocols in preventing unauthorized access to sensitive data.

Furthermore, the performance overhead introduced by cryptographic techniques needs to be considered. Evaluating the computational complexity and resource requirements of cryptographic operations in AI systems can help determine their feasibility and scalability in real-world applications.

Overall, evaluating cryptography for explainable AI requires a balanced assessment of transparency, privacy, and performance to ensure that cryptographic approaches effectively balance competing needs on accountability and privacy in data systems.

## VI.CONCLUSION

In conclusion, cryptography plays an important role in advancing the field of explainable AI by providing mechanisms to ensure transparency while maintaining data privacy and security. Using techniques such as homomorphic encryption, secure multi-party

computing (SMPC), zero-knowledge proofs (ZKP), and differential privacy, encryption methods enable AI models to operate act on sensitive data while ensuring privacy, security, and integrity. By encrypting both data and model parameters, the cryptography enables model inference and decision making on encrypted inputs without affecting the underlying information. This not only ensures that AI systems can provide explanations for their decisions, but also minimizes the risk of unauthorized access to sensitive data. Additionally, cryptographic protocols enable collaborative training of models on distributed data sources while preventing the exposure of individual data points, thereby promoting federated learning methods that ensure privacy protection. However, applying cryptographic techniques in explainable AI requires careful consideration of factors such as computational load, scalability, and the balance between transparency and performance. Despite these challenges, cryptography still offers a promising path to reconciling accountability and privacy requirements in AI systems, paving the way for the development of intelligent AI technologies. transparent and protect privacy. As research and development in this field continues to grow, cryptography will certainly remain the foundation of explainable AI, ensuring that AI systems can provide meaningful explanations while meeting the highest data security and privacy standards.

## REFERENCES

[1] dev Akeem Femi Kadri, Micheal Olaolu Arowolo, Ayisat Wuraola Yusuf-Asaju, Kafayat Odunayo Tajudeen, Kazeem Alagbe Gbolagade.A Roadmap for Enabling Industry 4.0 by Artificial Intelligence, 247-273, 2022

[2] Michael Moses Thiruthuvanathan, Kukatlapalli Pradeep Kumar, Noureen Nasar, KK Swathikiran, P Mani Joseph.XAI Based Intelligent Systems for Society 5.0, 365-400, 2024

[3] Dankan Gowda, Joohi Garg, Shaifali Garg, KDV Prasad, Sampathirao Suneetha,Innovative Machine Learning Applications for Cryptography, 27-46, 2024

[4] Andreas Holzinger, Anna Saranti, Christoph Molnar, Przemyslaw Biecek, Wojciech Samek.International workshop on extending explainable AI beyond deep models and classifiers, 13-38, 2022.

[5] Michael Moses Thiruthuvanathan, Kukatlapalli Pradeep Kumar, Noureen Nasar, KK Swathikiran, P Mani Joseph.XAI Based Intelligent Systems for Society 5.0, 365-400, 2024.

[6] Luca Vigano, Daniele Magazzeni.2020 IEEE European Symposium on Security and Privacy Workshops (Euros PW), 293-300, 2020.

[7] Andreas Holzinger 2018 world symposium on digital intelligence for systems and machines (DISA), 55-66, 2018.

[8] Dilip Kumar Sharma, Ningthoujam Chidananda Singh, Daneshwari A Noola, Amala Nirmal Doss, Janaki Sivakumar.Materials Today: Proceedings 51, 104-109, 2022.

[9] Samya Al Busafi, Basant Kumar,2020 9th International Conference System Modelling and Advancement in Research Trends (SMART), 323-327, 2020.

[10] Adrien Bennetot, Ivan Donadello, Ayoub El Qadi, Mauro Dragoni, Thomas Frossard, Benedikt Wagner, Anna Saranti, Silvia Tulli, Maria Trocan, Raja Chatila, Andreas Holzinger, A d'Avila Garcez, N Dıaz Rodrıguez.arXiv preprint arXiv:2111.14260, 2021.

[11] Akeem Femi Kadri, Micheal Olaolu Arowolo, Ayisat Wuraola Yusuf-Asaju, Kafayat Odunayo Tajudeen, Kazeem Alagbe Gbolagade.A Roadmap for Enabling Industry 4.0 by Artificial Intelligence, 247-273, 2022.