# A Survey on Trustworthy Edge Intelligence: From Security and Reliability To Transparency and Sustainability

Xiaojie Wang, Beibei Wang, Yu Wu, Zhaolong Ning, Song Guo, *Fellow, IEEE*, and Fei Richard Yu, *Fellow, IEEE*

*Abstract*—Edge Intelligence (EI) integrates Edge Computing (EC) and Artificial Intelligence (AI) to push the capabilities of AI to the network edge for real-time, efficient and secure intelligent decision-making and computation. However, EI faces various challenges due to resource constraints, heterogeneous network environments, and diverse service requirements of different applications, which together affect the trustworthiness of EI in the eyes of stakeholders. This survey comprehensively summarizes the characteristics, architecture, technologies, and solutions of trustworthy EI. Specifically, we first emphasize the need for trustworthy EI in the context of the trend toward large models. We then provide an initial definition of trustworthy EI, explore its key characteristics and give a multi-layered architecture for trustworthy EI. Then, we summarize several important issues that hinder the achievement of trustworthy EI. Subsequently, we present enabling technologies for trustworthy EI systems and provide an in-depth literature review of the state-of-the-art solutions for realizing the trustworthiness of EI. Finally, we discuss the corresponding research challenges and open issues.

*Index Terms*—Trustworthiness, edge computing, artificial intelligence, limited resources, interpretability.

## I. INTRODUCTION

With the development of Fifth-Generation (5G) wireless communication technologies, billions of wireless devices such as smartphones, sensors, and wearables can connect to the Internet. By 2025, the number of Internet of Things (IoT) devices are expected to reach 25.2 billion [1], and these devices will generate enormous data with more than 79 Zeta bytes per year [2]. Driven by the IoT, big data, and powerful computing, Artificial Intelligence (AI) has made further breakthroughs in intelligent applications such as natural language processing [3], computer vision [4], and robotics [5]. Especially in 2023, generative AI large models such as ChatGPT, DALL-E2, and DreamFusion, lead a trend that has become the focus of global attention. Generative large models based on billions of parameters can generate high-quality content in seconds, such as advertising images, short videos, writing copy, and voiceover. Currently, large models are mainly deployed on public clouds, mainly for training, with relatively small user sizes. However, the popularization of generative AI large models and the expansion of user scale will lead to a rapid increase in the computational demand for inference, exceeding the computational load for training.

In order to improve efficiency, reduce cost, and consider data security and privacy, Edge Computing (EC) [6] has emerged. Compared to cloud computing, EC deploys computational resources closer to users and data sources at the network edge, and thus enables low transmission latency and high communication efficiency. Edge Intelligence (EI) combines EC with AI techniques to fully leverage the advantages and potential of both. Enterprises, including Google, Microsoft, Intel and IBM, have developed pilot projects to demonstrate the benefits of EC in paving the last mile of AI [7]. For instance, MediaTek's APU 790 chip has a built-in hardware-level generative AI engine that supports high-speed and secure EI computation, and deeply adapts the transformer model for sub-acceleration, capable of generating images in less than one second. In addition, the hybrid precision quantization and memory compression techniques enable high-end smartphones to run Large Language Models (LLMs) with up to 33 billion parameters. While Google's release of PaLM2, a large model, realizes the full link of AI from the cloud to the network edge. Its lightweight version of Gecko runs on cell phones, significantly improving inference efficiency, reducing service costs, making the model applicable to a lot of application scenarios and users, and further promoting the development of EI deployment.

These efforts facilitate a wide range of AI applications, from real-time video analytics [8], Virtual Reality (VR), and Augmented Reality (AR) to smart healthcare [9], Autonomous Vehicles (AVs) [10] and Unmanned Aerial Vehicles (UAVs) [11]. EI encompasses the interconnection and collaboration of large-scale smart devices, sensors, and edge nodes so that data can be collected, processed, and exchanged across domains to support the digital intelligence transformation of various applications and businesses [12]. In this context, humans, objects and intelligences are interconnected to form a complex networked ecosystem.

However, with the computational load gradually migrate from the cloud to the network edge, issues of constrained bandwidth, storage and computational resources become prominent. This results in a series of challenges for EI ser-

X. Wang, B. Wang, Z. Ning (Corresponding author) are with School of Communications and Information Engineering, Chongqing University of Posts and Telecommunications, Chongqing 400065, China. E-mail: wangxj@cqupt.edu.cn, s220101143@stu.cqupt.edu.cn, ningzl@cqupt.edu.cn.

Y. Wu is with the School of Cyber Security and Information Law, Chongqing University of Posts and Telecommunications, Chongqing 400065, China. E-mail: wuy@cqupt.edu.cn.

S. Guo is with the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, Kowloon, Hong Kong, China. E-mail: songguo@cse.ust.hk.

F. R. Yu is with the Department of Systems and Computer Engineering, Carleton University, Ottawa, ON K1S 5B6, Canada. E-mail: RichardYu@cunet.carleton.ca.

TABLE I: Comparisons of related surveys on trustworthiness and EI.

| Topics | Ref. | Focuses | Contribution | Network scenarios |
|---|---|---|---|---|
| Trustworthiness | [13] | Requirements for AI trustworthiness | A survey on AI risk migration and system validation technologies to make AI trustworthy from the perspective of human involvement. | Not specific |
| | [14] | Technologies and applications for AI trustworthiness | A review of representative technologies and real-world application examples for trustworthy ML from a computational perspective. | Not specific |
| | [15] | Enhancing trustworthiness in AI product development process | A review of theoretical frameworks and systematic approaches to improve AI trustworthiness from the perspective of the entire life cycle of AI systems. | Not specific |
| EI | [16] | Efficient DL inference for edge devices | A review of state-of-the-art tools and techniques for efficient edge inference, highlighting the challenges and benefits of deploying DL models on resource-constrained edge devices. | Edge networks |
| | [17] | Compression techniques, hardware and software, and applications | A survey of challenges and solutions for deploying ML models on IoT devices with limited resources at the network edge. | Edge networks |
| | [18] | Algorithms for efficient distributed training | A summary of the challenges and techniques for communication-efficient EI, emphasizing collaboration between edge devices and servers to overcome communication overhead. | Edge networks |
| | [12] | The mutual support of advanced wireless technology and EI | A summary of solutions for EI in 6G networks when faced with challenges such as latency, energy consumption, network congestion, and privacy. | 6G networks |
| | [19] | Security | A summary of security threats and attack surfaces for IoT systems is presented and various ML and DL algorithms applied to IoT security are reviewed. | IoT |
| | [20] | Security and privacy | A survey of security and privacy issues related to network edge server deployment in 6G networks from the perspectives of EC, edge caching, and EI, respectively. | 6G networks |
| Trustworthy EI | This survey | Addressing security, reliability, transparency and sustainability of EI | A survey on definition, characteristics, architecture, technologies, solutions, challenges and, open issues for trustworthy EI from the perspective of the combination of AI and EC. | Edge networks |

vices, related to security, privacy, content generation accuracy, Quality of Service (QoS), and energy consumption [21]. These challenges directly impact the trust and acceptance of services by stakeholders, which range from end-users to technology developers, from government regulation to the general public. Therefore, research on trustworthy EI is of great practical significance to improve service quality, ensure system security, and promote sustainable development.

### A. Comparisons and Contributions

Although several studies have been devoted to providing solutions for trustworthy EI in recent years, there has not been a comprehensive review and synthesis of techniques and solutions for trustworthy EI.

As shown in Tab. I, some previous surveys focus on reviews of trustworthy AI [13]–[15]. Authors in [13] provide various requirements for ensuring the trustworthiness of AI, along with corresponding methods, all from a human-centered perspective. Authors in [14] present a detailed review of representative techniques for trustworthy AI from a computational point of view and discuss their practical applications in real-world scenarios. Differently, authors in [15] present a theoretical framework for important aspects of trustworthy AI from the

perspective of the entire life cycle of an AI system. Meanwhile, they systematically introduce available methods for realizing trustworthy AI. However, the above surveys primarily concentrate on plausibility in scenarios where computational resources are centralized. In contrast, our survey is centered around scenarios at the network edge that demand specialized responses to unique challenges, such as limited resources and network instability.

Several surveys provide reviews of EI architecture, technologies and optimization algorithms. For example, the work in [16] presents a survey of software tools for hardware-algorithm co-design, along with an analysis of existing choices and trade-offs in hardware platforms. Authors in [17] delve into a range of applications related to Machine Learning (ML) and EC from an operational perspective. Furthermore, authors in [18] summarize communication efficient algorithms for distributed training of AI models. Authors in [12] focus on EI in 6G networks.

Differently, researchers in [19] and [20] review security and privacy protection strategies at the network edge. To be specific, authors in [19] provide a review of the state-of-the-art ML and Deep Learning (DL) approaches from the perspective of IoT security. Authors in [20] explore security threats and
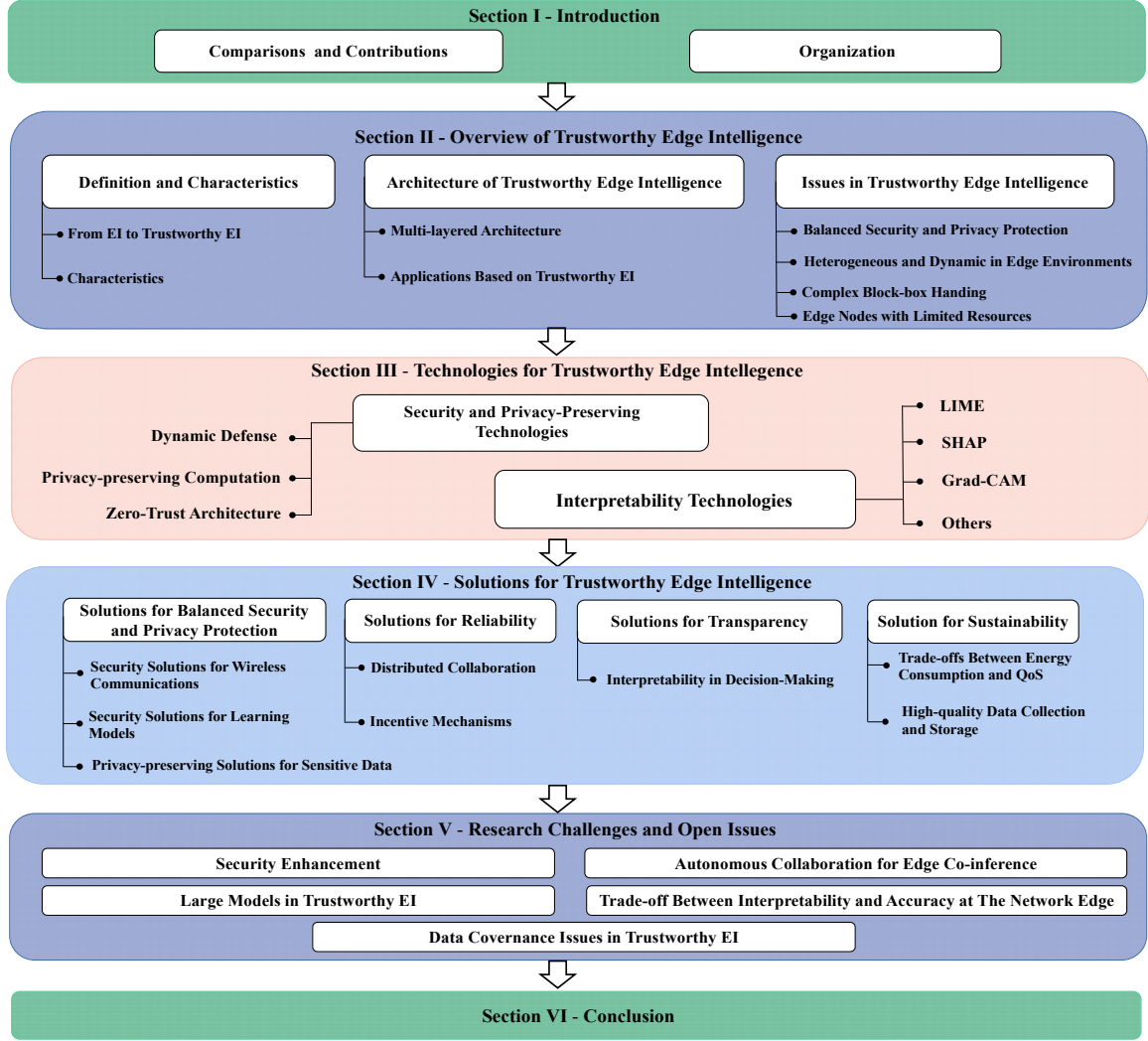
Fig. 1: Structure of the survey.

countermeasures associated with EC, edge caching, and EI in the context of 6G network edges.

In summary, existing surveys focus on solving edge-efficient inference and training algorithms as well as security and privacy challenges. However, there is a lack of systematic survey that provide a comprehensive and professional discussion of the concepts, essential features, techniques, solutions, and challenges of trustworthy EI. ***To the best of our knowledge, this survey is the first to provide a comprehensive summary of trustworthy EI from a combined EC and AI perspective.*** The contributions of this survey can be summarized as follows:

- We first provide a definition of trustworthy EI, its essential characteristics, and establish a three-layer architecture to support the concept. We then summarize some key issues with respect to the vulnerability of edge networks and AI models to highlight the difficulties faced in achieving trustworthiness of EI.
- We introduce enabling techniques for trustworthy EI in terms of security and interpretability, respectively.
- We provide a comprehensive and in-depth investigation of recent studies on trustworthy EI based on its issues

and requirements. In addition, lessons learned for each kind of approaches are also provided.
- Finally, we provide a detailed discussion about our vision for the future development of trustworthy EI, and identify a number of open challenges that may give rise to promising research directions.

### B. Organization

As shown in Fig. 1, the survey is organized as follows. Section II examines the concept, characteristics, architecture, and issues of trustworthy EI. We present technologies used to achieve trustworthiness in Section III, and discuss solutions to realize trustworthy EI in Section IV. Challenges and open issues for trustworthy EI are provided in Section V, and the survey is concluded in Section VI.

## II. OVERVIEW OF TRUSTWORTHY EDGE INTELLIGENCE

This section provides a comprehensive overview of trustworthy EI, covering its definition, characteristic, architecture, cutting-edge application scenarios and issues.

## A. Definition and Characteristics

In this following, we begin by examining the problems faced by EI, clarifying the urgency of trustworthy research. Subsequently, we explicitly define the concept of trustworthy EI and the fundamental characteristics it possesses.

*1) From EI to Trustworthy EI:* In the actual deployment and application of EI, users and organizations remain skeptical and cautious as it introduces several concerns and risks: i) Security and privacy concern: In open edge environments, communication vulnerabilities may lead to sensitive data leakage, while AI algorithms are susceptible to poisoning and adversarial attacks, especially large models that are pre-trained and fine-tuned at the network edge, increasing the risk of poisoning and adversarial attacks. As a result, users may lose confidence in EI due to concerns about the unreliability of models in mission-critical situations, such as parking and acceleration decisions for AVs; ii) Imbalanced performance: Limited computational resources, restricted storage capacity, and limited energy supply of edge devices lead to imbalance in the performance of EI services, especially in real-time applications. Because reducing the model complexity to adapt to resource constraints may bring about degradation in model accuracy and may even increase the risk of security and legal liabilities; iii) Un-interpretable model: Black-box models, such as LLMs usually set a large number of parameters to boost inference performance, making it difficult for users to understand how the model works, its limitations, and potential flaws. Even open source tools like LLaMA provide only limited interpretability. Interpreting and publicizing training data becomes very complicated, especially for proprietary models such as ChatGPT and Claude, whose architectures and training data are not yet publicly available [22].

These issues restrict the widespread application of EI and introduce a trust deficit for the digital network ecosystem. Therefore, it is important to investigate trustworthy EI to alleviate the trust concerns. Before defining trustworthy EI, we draw on the insights from [23] to briefly discuss the concepts, distinctions, and connections between trust and trustworthiness. Although their precise definitions depend on specific application contexts, there is a general and broad consensus that trust is a psychologically and emotionally dimensional concept involving an individual's or organization's confidence in and expectations of another person, organization, or system. Trust is formed when an individual or organization feels reliable about a party's behavior, commitment, or competence. In contrast, trustworthiness focuses on the objective level and relies on the reliability, stability, and security of a system, service, or individual in its design, implementation, and operation. A trustworthy entity should be able to fulfill its commitments, protect the interests of others, and demonstrate reliability and accountability in all aspects. Thus, judgments of trust reflect beliefs about the trustworthiness of another party. Ensuring that the entity has a high degree of trustworthiness is the basis for building trust among users and other stakeholders.

In the context of EI, this entity represents systems that integrate EC, IoT devices, edge servers, and AI models deployed at the network edge. It aims to enable real-time intelligent decision-making and computation locally or in close proximity to the data source. Therefore, trustworthy EI can be defined by *"A EI system employs a variety of technologies and strategies to ensure the system to achieve the predefined goals during the design phase while minimizing potential risks and hazards"*.

*2) Characteristics:* While considering the unique features of EI, such as limited resources, real-time requirements, and distributed deployment, and simultaneously conducting a comprehensive evaluation of the overall performance of EI systems, the characteristics of trustworthy EI can be briefly summarized as follows: Endogenous security, reliability, transparency, and sustainability. These characteristics not only address current issues faced in EI systems but also effectively respond to potential new challenges in the future, enabling the system to possess enhanced adaptability and garner trust from stakeholders. We elaborate these characteristics in the following content:

- **Endogenous Security:** In the field of EI, traditional security approaches face the dilemma of coping with the challenges of real-time, distributed, and limited resources. Endogenous security emphasizes that security mechanisms should have the ability to proactively adapt to dynamic changes at the network edge with lightweight deployment. By introducing technologies and policies such as Physical Layer Security (PLS), zero trust, and anomaly detection, trustworthy EI systems are able to provide cost-effective, efficient, and reliable security in areas such as wireless communications, distributed learning, and data sharing. The endogenous security breaks through the limitations of traditional security methods, emphasizes the inherent initiative and intelligence of the system, and brings a new perspective to the security of trustworthy EI systems.

- **Reliability:** It typically involves the system's fault-tolerance capability and performance consistency. Different from traditional networks, trustworthy EI systems place special emphasis on real-time requirements, resource constraints, and distributed computing. Firstly, many EI-driven applications such as AVs, AR, and VR demand high real-time performance and low latency [24]. The trustworthy EI must ensure real-time inference for quick decision-making within extremely short timeframes. This not only depends on communication methods and algorithm complexity but also relies on efficient resource scheduling on edge devices. In contrast, core networks usually operate in large data centers with relatively abundant computing and storage resources, may prefer the reliability of communication links. Lastly, reliability needs to account for the collaborative execution of tasks across different edge nodes to ensure the overall stability of the system.

- **Transparency:** In trustworthy EI system, transparency encompasses not only the interpretability of decision-making process but also clear data presentation. First, it is crucial to ensure that the decision-making process is interpretable. Given the limited resources of edge devices, interpretable techniques and strategies must strike a bal-
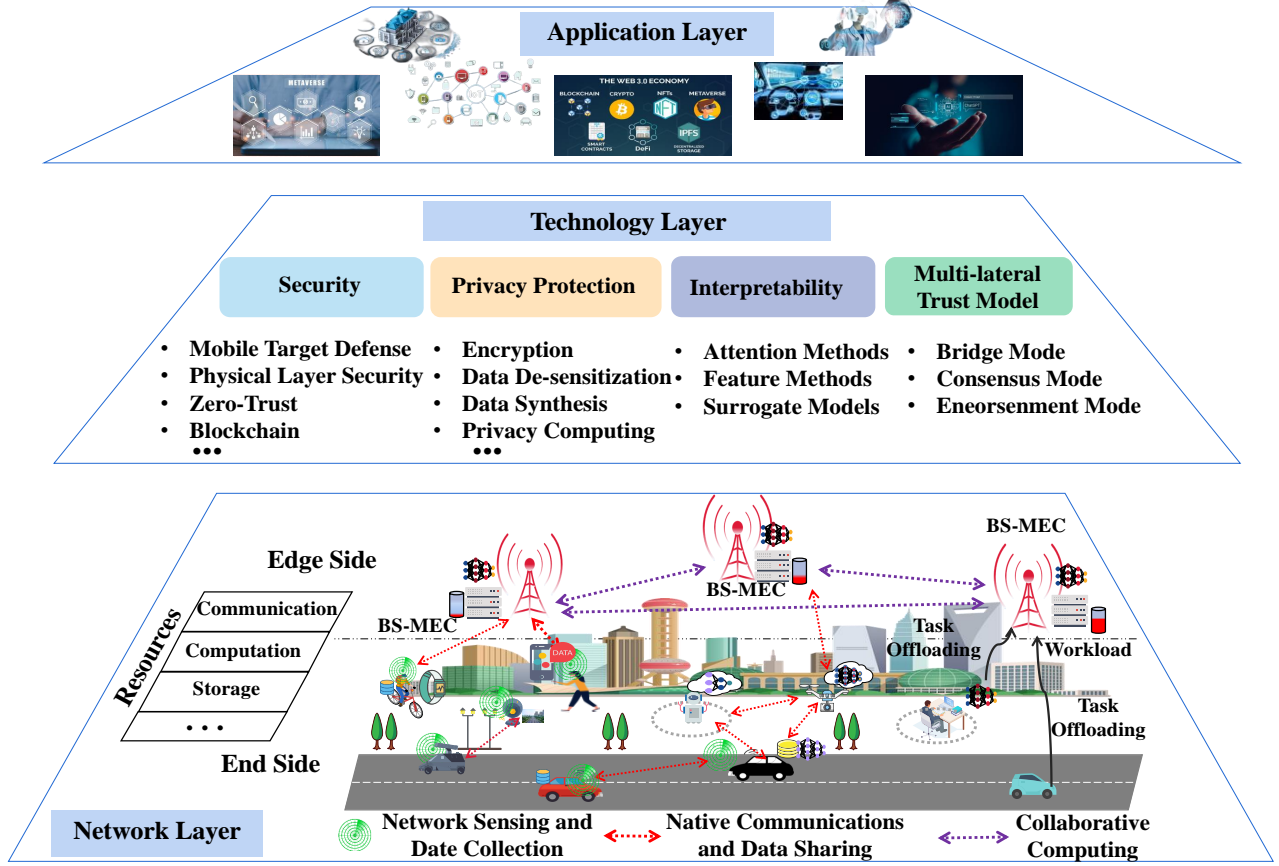
Fig. 2: Architecture of trustworthy EI: The network layer integrates communication, computation, perception, and intelligence, delivering the application layer high-performance, interpretable decision-making, resilient and sustainable solutions. Meanwhile, the technology layer establishes a trustworthy environment, offering robust technical support across the network and application layers, fostering seamless collaboration for comprehensive trustworthy EI.

ance, ensuring that explanatory information is both comprehensive and adheres to computational constraints [25]. Second, transparency in data usage allows users to understand key information about the data source, quality, and processing within the system. This transparency is essential for establishing user trust, ensuring privacy compliance, and enhancing system credibility.

- **Sustainability:** It reflects the long-term health and stability of the trustworthy EI system. On the one hand, edge devices rely on limited energy sources such as batteries, and the use of energy-efficient algorithms can extend device lifespan and reduce environmental stress [26]. In addition, the distributed infrastructure of EI systems includes outdoor environments where adverse environmental impacts can be reduced through low-power hardware and green computing practices. On the other hand, EI utilizes edge devices to continuously sense and collect data, emphasizing the importance of high-quality data resources to ensure model accuracy and reliability [27]. In conclusion, sustainability not only promotes environmental friendliness, but also enables trustworthy EI systems to provide long-lasting and reliable services while maintaining peak performance.

In summary, the above characteristics constitute the trustworthiness of EI, fostering trust between the EI system and stakeholders. They are usually interrelated rather than independent. In resource-limited edge networks, security measures should be lightweight. Reliability allows the system to balance performance and safety, while security enhances resistance against external threats, improving overall reliability and transparency. In addition, transparency in decision-making empowers users and developers to understand how the model works and facilitates prompt adjustments in case of errors. A secure, reliable, and transparent system is more likely to sustain in a dynamic environment, showcasing mutual reinforcements and trade-offs among these characteristics for trustworthy EI.

### B. Architecture of Trustworthy Edge Intelligence

In this subsection, a multi-layered architecture for scalable and trustworthy EI systems is presented. Following, we detail the components of this architecture and the supported preamble application services.

*1) Multi-layered Architecture:* As shown in Fig. 2, the architecture is divided into the network layer, the technology layer and the application layer.

The network layer, as a basic and key level, covers edge and IoT devices. Its main task lies in realizing cooperative communication among devices and accomplishing data collection, transmission, processing and analysis at the network edge, so as to build the foundation of EI service [2]. Among them, end devices with embedded sensing, communication and limited processing capabilities collect data in the environment, establish connections, exchange data and share models with other devices at the network edge. Edge nodes include fog nodes, gateways, and EC servers. Compared to end-side devices, edge nodes have more powerful storage and computation capabilities and can provide high-quality network and processing services with low latency. The edge-side ones are mainly responsible for end-side and edge-side resource fusion, coordinating heterogeneous resources to achieve fast and flexible service deployment and low-latency task processing.

The technology layer plays an integral role in trustworthy EI, providing critical trustworthy technology support to the network layer. By ensuring high-performance services, sustainable and interpretable intelligent decision-making, and a reliable and resilient network, the technology layer safeguards the trustworthiness and security of the entire system. In particular, the multi-lateral trust model plays a pivotal role in establishing a robust foundation of trust for communication, collaboration, and information sharing within the network layer. It achieves this by creating and managing a framework for trust relationships among devices. The model encompasses three key modes: bridge, endorsement, and consensus, as shown in Fig. 3.

In the "bridge" model, a central authorization authority conducts peer-to-peer authentication between entities A and B, facilitating the direct establishment of a trust relationship. The "endorsement" mode relies on a third-party organization to assess the trustworthiness of an entity, subsequently transmitting the evaluation results of entity A to entity B to establish trust. The "consensus" mode is considered as the most crucial one, employing distributed transactions among entities. This model offers a flexible and scalable solution for building trust relationships in a distributed manner.

At last, the application layer makes full use of the trustworthy technologies and infrastructure provided by the network layer to realize high-quality, well-accepted and trustworthy application services. In the following subsection, we elaborate on several applications in detail, demonstrating the benefits of trustworthy EI in these emerging scenarios.

*2) Applications based on Trustworthy EI:* Trustworthy EI is at the forefront of innovative applications in multiple cutting-edge fields. From mobile Artificial Intelligence-Generated Content (AIGC) to Web 3.0 and onward to Metaverse, these domains are continuously expanding the boundaries of our digitized lives. In the following, we delve into these trendsetting application areas and analyze the key role of trustworthy EI in ensuring the secure, efficient, and reliable deployment of mobile AIGC, Web 3.0, and the Metaverse.

- **Mobile AIGC:** AIGC has become a novel approach for processing and manipulating data, supporting multi-modal input and output to automatically generate creative and personalized content. The mobile AIGC is the
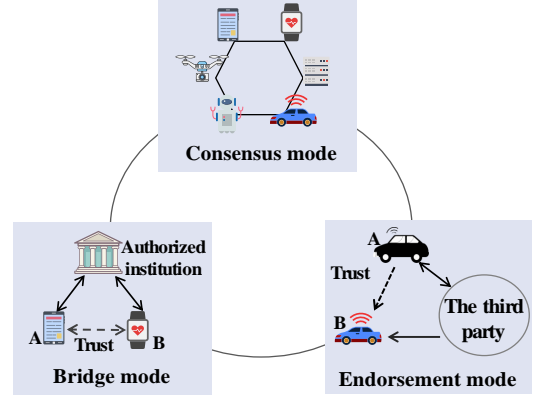


Fig. 3: The multi-lateral trust model.

deployment of AIGC services in mobile edge networks to enhance the user experience of mobile applications and generate innovative solutions [28]. The trustworthy EI architecture plays a crucial role to ensure efficient, accurate, and secure mobile AIGC services. Specifically, in the trustworthy EI, communication and distributed computational resources at the network edge and the end side are integrated and allocated in real time, enabling mobile AIGC applications to quickly adapt to user demands, maximize the utilization ratio of limited resources, and maintain accurate content generation. Lightweight security and privacy protection mechanisms ensure that AI model training is safe while minimizing the use of computational resources, thus contributing to improving efficiency, security, and accuracy of mobile AIGC services. Therefore, the reliability and security characteristics of trustworthy EI strongly support the secure and sustainable deployment of mobile AIGC.

- **Web 3.0:** It is considered that the next generation of the Internet allows users to read, write, and own content [29]. Current Web 3.0 services primarily focus on supporting blockchain applications, including Decentralized Identifiers (DIDs), digital asset management, Decentralized Applications (DApps), cryptocurrency-based Decentralized Finance services (DeFi), and Metaverse. Trustworthy EI plays a crucial role in realizing the decentralized, secure, and user-controlled Web 3.0. First, it allows data processing on edge devices, providing Web 3.0 users with a secure means to own and control their digital identity and data. Second, decentralized applications in Web 3.0 may involve edge devices such as IoT devices and sensors. The trustworthy EI enables these devices to efficiently and reliably collaborate in executing smart contracts and processing local data, thereby reducing network latency. In summary, trustworthy EI provides essential technical support for the implementation of Web 3.0, contributing to the creation of a decentralized, user-friendly, and secure digital ecosystem.

- **Metaverse:** It is considered as the evolutionary paradigm of Web 3.0. This concept integrates a plethora of existing technologies, including 5G, Multi-access Edge Computing (MEC), AI, VR, blockchain, digital currency, IoT and

human-computer interaction [30]. Trustworthy EI significantly enhances the computational efficiency, accuracy, scalability, privacy, and security of AI-driven virtual services in the Metaverse, including AR/VR recommendation and cognitive virtual identities. Specifically, trustworthy EI supports distributed collaborative computational and real-time resource allocation, markedly improving the processing efficiency of tasks such as high-dimensional data processing, 3D virtual world rendering, and avatar computation in the Metaverse. Edge learning with secure aggregation mechanisms facilitates local training, effectively reducing communication costs associated with AR content delivery, such as 3D objects or high-resolution video streams, and ensuring the security of model training and privacy of sensitive data. Furthermore, trustworthy EI with transparency helps users to understand why the system makes a particular recommendation or behavior, thus increasing their trust and acceptance of the virtual service. In summary, the trustworthy EI is crucial for ensuring QoS, user privacy, security, and stable operation of Metaverse applications.

### C. Issues in Trustworthy Edge Intelligence

Although trustworthy EI can provide many benefits, in the face of resource-constrained edge networks and complex AI models, in order to realize the trustworthiness of EI, the following issues need to be considered.

*1) Balanced Security and Privacy Protection:* The escalating demand from users for real-time and seamless services and requests necessitates that trustworthy EI systems strike a balance between security and privacy protection capabilities and user experience [20]. Excessively strict security policies may result in resource wastage, reduced system flexibility, and even negative impacts on user privacy. By appropriately configuring network architecture and employing intelligent security mechanisms, the system can effectively resist potential security risks while maintaining a high level of network and service quality. Hence, a flexible and balanced security strategy becomes particularly crucial in trustworthy EI environments. We discuss solutions to solve this issue in subsection IV-A.

*2) Heterogeneity and Dynamics in Edge Environments:* The differences in computing capability, storage capacities, network bandwidth, sensing abilities, and domain knowledge of each node make the management and scheduling of resources extremely complex. Therefore, it is necessary to further consider how to coordinate and fully utilize different resources and knowledge of nodes to improve the performance and effectiveness of the whole system. In addition, the resource state of nodes changes over time, including device connection and battery power, further increases the difficulty of resource management. This requires the system to be able to flexibly respond to the dynamic changes of nodes to ensure the normal operation at different time and conditions. We thoroughly examine solutions to this issue in subsection IV-B.

*3) Complex Black-box Handling:* According to the previous discussion, to enhance the transparency of black-box models, interpretable methods need to be provided. However, increasing interpretability often introduce additional computational and storage costs, which can be a challenge for the limited resources of edge devices. Interpretability methods for complex black-box models may entail sacrificing real-time performance. Therefore, it is important to balance the interpretability requirements and real-time performance. We delve into detailed solutions for this issue and provide a summary of relevant literature in subsection IV-C.

*4) Edge Nodes with Limited Resources:* Edge devices are often constrained by limited computation and storage resources, and executing complex AI inference and training tasks may consume a significant amount of these resources, thereby reducing the lifetime of these devices. Specific EI applications, such as AVs and smart factories, have extremely high requirements for latency, which may result in high power consumption, adversely affecting energy efficiency. In addition, certain intensive AI tasks can only be performed on cloud servers, which causes high communication costs. Thus, leveraging ubiquitous network resources while reducing the communication overhead becomes a major bottleneck for AI applications at the network edge. We discuss some possible solutions for this issue in subsection IV-D.

## III. TECHNOLOGIES FOR TRUSTWORTHY EDGE INTELLIGENCE

In this section, we focus on technologies that play a crucial role in ensuring the trustworthiness of EI. The utilization of these technologies serves as a means to address existing challenges and enhance system trustworthiness.

### A. Security and Privacy-Preserving Technologies

The devices involved in EI are usually deployed close to the user side to realize real-time data processing and decision making. Both the security of computation and the risk of sensitive data leakage may reduce user trust. Therefore, security and privacy-preserving technologies can be utilized to minimize the susceptibility of EI systems to attacks and prevent data leakage, thereby enhancing user confidence in the system. Next, we provide a detailed description of security and privacy-preserving techniques, respectively.

*1) Dynamic Defense:* With the automation and intelligence of adversarial techniques continually advancing, dynamic defense technologies gain widespread attention among researchers. The idea is to dynamically alter and disguise characteristics of target systems, thereby increasing attack costs, enhancing the system resilience, and bolstering defense capabilities. In the following, we introduce these related defense technologies in detail, which include Cyberspace Mimic Defense (CMD), Mobile Target Defense (MTD), and Cyber Deception (CD).

- *Cyberspace Mimic Defense:* The CMD aims to prevent attackers from forming effective attacks by means of conditional evasion, so that the inevitable endogenous security problems do not become security threats to the system. The core idea is to organize multiple redundant and heterogeneous methods to jointly process external

requests, and achieve dynamic scheduling through negative feedback to compensate for security flaws in the cyberspace.

- *Mobile Target Defense:* The MTD is an active defense mechanism that prevents network attacks by continuously and dynamically changing attack surface [31]. Its objective is to create uncertainty for attackers and shift the asymmetry between attackers and defenders.
- *Cyber Deception:* Compared to MTD, CD employ more aggressive strategies, intentionally providing false information (such as baits and honeypots) to mislead attackers [32].

*2) Privacy-preserving Computation:* Ensuring data privacy usually requires the use of different techniques and methods, such as Differential Privacy (DP), Homomorphic Encryption (HE), Secure Multi-party Computation (SMC), and Trusted Execution Environment (TEE). In the following, we provide a brief introduction to them.

- *Differential Privacy:* DP aims to provide statistical guarantees for individual data while minimizing the disclosure of individual privacy [33]. The main idea is to add noise, such as Laplace noise, to the original query results (numerical or discrete values). The added noise prevents the inference of significant information about individuals from query results, preserving personal privacy.
- *Homomorphic Encryption:* HE is a class of encryption mechanisms that support processing and computation of ciphertexts [34]. In actual training and inference process, HE techniques can ensure the security of model parameters and raw data, thus developing trustworthy EI models [35]. However, existing HE solutions need to address the problem of high computational overhead, especially in resource-constrained edge environments.
- *Secure Multi-party Computation:* SMC enables collaborative computation on a combined dataset without compromising the data privacy of individual parties [36]. By implementing SMC at the edge nodes, trustworthy collaborative computation can be achieved, and data availability without data visibility can be guaranteed.
- *Trusted Execution Environment:* It is an isolated processing environment designed to provide computing and storage capabilities with security and integrity guarantees. The basic idea is to allocate segregated hardware memory for sensitive data, ensuring secure transmission, storage, and processing. Therefore, deploying TEE can ensure the confidentiality of both EI models and input data [37].

*3) Zero-Trust Architecture:* The fundamental principle of zero-trust is to distrust any access, whether inside or outside the network, and to advocate for continuous authentication and dynamic authorization to ensure global defense. Zero-trust network access relies on micro-segmentation and network isolation, eliminating the need for a virtual private network by granting access to the network only after thorough verification and authentication.

The increased connectivity at the network edge inevitably expands attack surface, allowing attackers to access systems through multiple entry points. Additionally, not all devices can receive security updates promptly, potentially enabling attackers to exploit multiple vulnerabilities to gain access to the network. Therefore, applying the zero-trust principle to edge networks forms a new defense boundary, capable of integrating security and networking anywhere [38]. Zero-trust edge, based on continuous verification of user identity and context, provides explicit access to applications.

*4) Others: Physical Layer Security (PLS)* is used to enhance the security of wireless communication systems. Different from traditional encryption methods, PLS relies on physical properties of communication channels rather than algorithms and keys. Common PLS techniques include: artificial noise [39], cooperative jamming [40], and beamforming [41]. Specifically, artificial noise is used to mask original data from eavesdroppers by intentionally adding noise to the transmitted signal in the communication channel. In cooperative jamming, multiple nodes work together to protect communication privacy by interfering with potential eavesdroppers. Beamforming technology concentrates signal energy in a specific direction while reducing signal strength in other directions, making it difficult for eavesdroppers to intercept.

*Anomaly detection* is used to automatically identify abnormal data. During the process of edge training, data updates from various training nodes can be analyzed to detect malicious updates based on differences between pairs of remote updates [42]. This process enables the acquisition of a robust global model.

*Blockchain* is a cryptographic, decentralized, and user-transparent technology that provides secure transactions and computing at the network edge [43]. Compared to the above techniques, blockchain has unique technical features such as consensus protocols and distributed ledgers. These features enable blockchain to effectively regulate security risks in EI systems. First, the consensus mechanism ensures the establishment of trust among devices for model training. Second, the tamper-proof distributed ledger keeps the recording of authentic and reliable information, promoting a transparent process [44]. Last, it rewards participating nodes based on their contributions, which incentivizes selfish nodes to provide their local resources [45].

### B. Interpretability Technologies

Local Interpretable Model-agnostic Explanations (LIME) and SHapley Additive exPlanations (SHAP) are applicable to various models, while Gradient-Weighted Class Activation Mapping (Grad-CAM) specifically targets Convolutional Neural Network (CNN) models. In the following, we offer a brief overview of these technologies.

*1) LIME:* It is a model-agnostic interpretability method, and focuses on interpreting individual instances rather than the entire model to provide local explanations for specific test inputs [46]. It introduces random perturbations to the input instances of a black-box model and trains an interpretable surrogate model, such as decision trees and linear models. The weights of the surrogate model can directly reflect the significance of features and their impacts [47]. To ensure interpretability and local fidelity, LIME tries to minimize the

discrepancy between the surrogate model and the black-box model at the instance point.

*2) SHAP:* It can be regarded as a unified approach that combines LIME and shapley values [48]. The shapley value of a single feature is the weighted average of the marginal contribution of that feature to a subset of all feature combinations. SHAP is the basis for fairly distributing contributions of each feature to the model and has three desirable properties (i.e., local accuracy, missingness, and consistency) [49]. Generally, SHAP is applied as an interpretability method based on feature correlation in edge scenarios such as smart healthcare [50] and intrusion detection [51], not only to interpret final decisions, but also to support industry experts to quickly optimize and evaluate the correctness of their judgments.

*3) Grad-CAM:* It is a visual local interpretability method designed specifically for CNN models [52]. It calculates the gradient of the target class with respect to the last convolutional feature map of the CNN. This gradient information helps to identify image regions that have the biggest contribution to model predictions. By back-propagating gradients and multiplying them with the feature map, the importance weight is assigned to each pixel. These weights are then used to generate a heat map, which highlights the significant regions in the image.

Grad-CAM supports a wide range of CNN models and does not require further changes to the model architecture. Since the size of feature maps is usually much smaller than the input image, heat maps produced by Grad-CAM may not provide precise localization. Authors in [52] combine a fine-grained visualization method of guided back-propagation with Grad-CAM to produce high-resolution activation maps. The guided Grad-CAM tries to provide intuitive and effective interpretation in clinical medical image analysis, helping physicians to determine the location, type, and severity of lesions, and facilitating the advancement and application in the field of medical imaging analysis [53].

*4) Others: Attention mechanism* is a widely used technique that mimics features in the human visual and perceptual system to process and interpret input data. Visual attention is a method for visualizing the attention weights of a model, usually for text and image data.

*Rule-based* interpretability techniques explain the behavior of a model by defining a set of rules. These rules can be created manually or generated by automatic learning techniques. Representative techniques include decision trees, rule sets, and expert systems.

## IV. SOLUTIONS FOR TRUSTWORTHY EDGE INTELLIGENCE

The realization of trustworthy EI needs to satisfy the four key characteristics of security, reliability, transparency, and sustainability. In order to ensure these trustworthy characterizations in complex edge environments, it is necessary to address the issues described in subsection II-C. Thus, this section aims to provide an in-depth discussion of existing approaches.

### A. Solutions for Balanced Security and Privacy Protection

In this subsection, we focus on discussing solutions for the security issues mentioned in subsection III-C-1 to establish endogenous security capabilities for trustworthy EI systems, as show in Fig. 4. First, we will investigate security solutions for wireless communications to ensure secure and reliable services in dynamic environments. Second, we will delve into security solutions for learning models to guarantee the integrity and availability of model training and inference in distributed environments. Finally, we will explore privacy-preserving solutions for sensitive data. These solutions aim to strike a balance between security and user experience, avoiding wasted resources and reduced system flexibility.

*1) Security Solutions for Wireless Communications:* Researchers usually adopt strategies such as enhanced encryption protocols and secure communication channels to ensure the confidentiality and integrity of data transmission. We summarize the related studies in Tab. II.

**Lightweight encryption:** Traditional encryption strategies, such as HE, may not be applicable to resource-constrained edge devices due to their computational requirements [20]. Authors in [54] enhance the security of EI by developing a lightweight and leak-resistant certified key exchange protocol. Meanwhile, the proposed protocol is designed for mainstreamed communication standards. Authors in [55] propose an Attribute-Based Multi-Authority Signed Confidentiality (ABSC) scheme for IoT data sharing. This scheme offloads most of the computation to edge servers, reduces communication and storage costs through short and constant-size ciphertexts, and employs a hierarchical multi-privilege architecture. Unlike the above studies, authors in [56] utilize Quantum Key Distribution (QKD) to protect confidential data for Semantic Information Communication (SIC). In order to reduce the cost of edge users, they use a stochastic planning model to optimize resource allocation and use Shapley values to distribute the cost among QKD service providers.

**Secure communication channels:** PLS is a classic method, unlike encryption, it does not face the challenges of complex key management and high-overhead key distribution. However, considering constraints such as latency, power consumption, and QoS, security policies based on PLS require trade-offs between system performance and security capabilities.

Authors in [39] utilize devices that do not participate in Federated Learning (FL), such as Sensor Nodes (SNs), to send jamming signals to defend against eavesdropping attacks. They optimize local training time, model upload time, and transmission power to obtain the best pairing of training nodes and SNs. Considering that friendly jammers have limited power, eavesdroppers are not able to eavesdrop on all channels. Therefore, authors in [40] propose a cooperative jamming power allocation strategy based on game theory that considers eavesdroppers as strategic players.

Unlike PLS, Intelligent Reflective Surface (IRS) technology can support high-dimensional data transmission while enhancing the security of wireless communication networks by adjusting signal amplitudes and phases. Authors in [41] propose a Deep Reinforcement Learning (DRL)-based approach to jointly optimize the beam-forming matrix at the base station
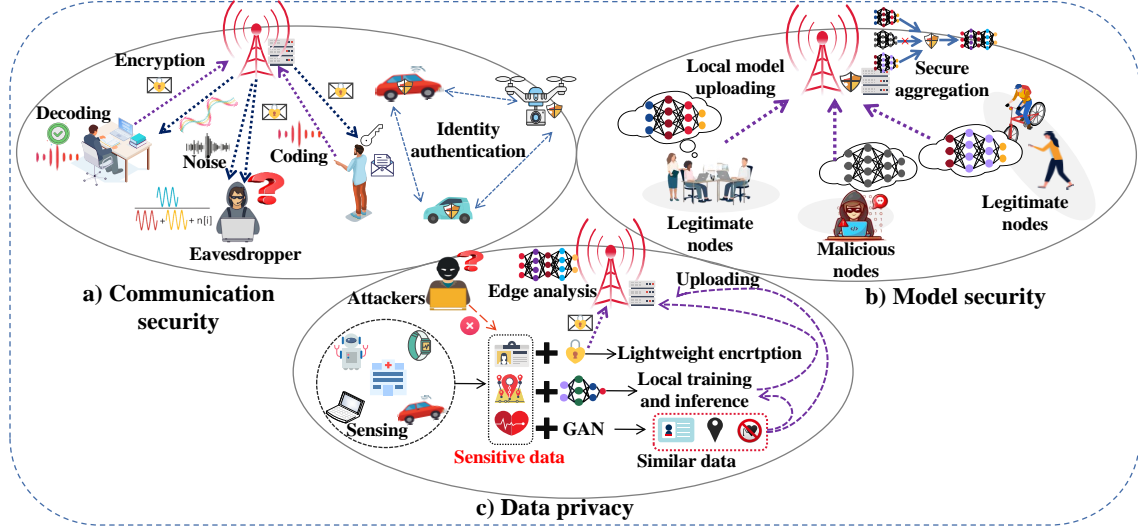
Fig. 4: Illustrative solutions of balanced security and privacy protection: a) Security solutions for wireless communications ensure availability of the edge network and security of data transmission; b) Security solutions for learning models address threats such as model poisoning and adversary attacks, ensuring trustworthiness and robustness of learning models; and c) Privacy-preserving solutions focus on handling sensitive data at the network edge.

and the reflected beam-forming matrix at the IRS. It aims to maximize secure communication to legitimate users in a dynamic environment while meeting QoS requirements. With the advancement of technology, some intelligent attackers are able to selectively employ eavesdropping or active jamming attacks [57]. Thus, authors in [57] consider a non-cooperative game between a base station and an intelligent attacker. They jointly optimize power allocation and beamforming to improve the secrecy rate. Meanwhile, RL is used to predict the attack behavior and selectively send artificial noise signals.

Differently, authors in [58] utilize cyber deception to send true information to the intended receiver while injecting fake information to confuse eavesdroppers. Traps are strategically deployed to attract eavesdroppers and provide them with increasingly clear fake messages, thus establishing a secure communication channel between senders and receivers. This method ensures security of exchanged information, even if eavesdroppers gain access to the secret channel information.

**Identity authentication:** Due to the increase in the number of wireless users and the openness of electromagnetic wave propagation, wireless communications are vulnerable to security threats. Therefore, ensuring the legitimacy of identities of communicating parties is essential for preventing various attacks and securing the content of communications [60].

To accommodate the high dynamics and resource constraints of edge networks, collaborative approaches can provide accurate and robust authentication by sharing multi-dimensional information among devices. For example, authors in [59] propose a collaborative physical layer authentication method based on FL. The scheme utilizes a set of reputable edge devices to co-construct the authenticator to eliminate the limitations of insufficient resources and high authentication complexity. At the same time, the scheme collects multi-dimensional channel state information from different collab-

orating devices, making it difficult for external adversaries to simulate, thus improving the endogenous security of the system. Similarly, the work in [60] assists service providers in verifying user identities by collecting and processing user location-related characteristics, such as received signal strength and movement trajectories. This article also proposes a context-aware group updating algorithm for adaptively updating cooperative peers and authentication features in dynamic networks.

Blockchain-based authentication is also a new way of identity authentication [63]. Authors in [61] design an efficient blockchain-based anonymous cross-domain authentication scheme to enable reliable communications among cross-domain IoT devices. By combining blockchain and dynamic accumulator technology, the scheme achieves fast authentication, reducing computational pressure for edge devices while ensuring device anonymity. Similarly, authors in [62] propose a complete cross-domain authentication and privacy protection scheme based on federated blockchain. The process of cross-domain authentication is divided into three phases: authorization, on-demand pseudo-identity generation, and identity authentication. In the intra-domain authentication phase, it performs according to the original scheme adopted by the domain. This not only provides excellent scalability, but also greatly reduces the deployment cost and time overhead of the system.

Unlike the above studies, the zero-trust framework aims to provide continuous authentication. Authors in [38] propose a zero-trust and EI-based scheme. Proactive and continuous authentication is achieved by periodically monitoring and re-evaluating variable attributes throughout the request lifecycle. In addition, the article employs an edge intelligence algorithm based on neural-supported decision trees to improve authentication accuracy.

TABLE II: Security solutions for wireless communications.

| Solutions | Ref. | Description | Optimization Metrics |
|---|---|---|---|
| Lightweight encryption | [54] | A lightweight authentication key agreement to enhance EI security. | Security and computation costs |
| | [55] | A lightweight ABSC scheme to ensure secure data sharing in IoT. | Security and computation costs |
| | [56] | A two-stage stochastic optimization model for low-cost QKD-SIC. | Security, computation and commutation costs |
| Secure communication channels | [39] | A channel sharing scheme to improve secure FL by incentivizing idle devices to send jamming signals. | Security and FL latency |
| | [40] | An iterative power allocation strategy to defend against eavesdroppers. | Security |
| | [41] | A DRL-based beamforming method for IRS-assisted secure communication. | Security and QoS |
| | [57] | A DQN-based approach that combines IRS and RL to counter intelligent attackers. | Security |
| | [58] | An active defense method that uploads confidential information while sending "fake" but meaningful information to confuse the eavesdropper. | Security |
| Identity authentication | [59] | A collaborative FL-based authentication method for endogenous system security. | Security, privacy, and commutation costs |
| | [60] | A new cooperative authentication scheme helps service providers to authenticate subscribers in a distributed manner. | Security and latency |
| | [61] | An efficient blockchain-based anonymous cross-domain authentication scheme for reliable communications among cross-domain IoT devices. | Security and commutation costs |
| | [62] | A blockchain-based secure and privacy-preserving authentication scheme. | Security, privacy, and commutation costs |
| | [38] | A zero-trust and EI based continuous authentication method for satellite networks. | Security and accuracy |

*2) Security Solutions for Learning Models:* In EI, poisoning attacks and adversarial attacks are prevalent methods of attack. In response to these threats, we provide an overview of the primary defense methods against each of these attack types. In Tab. III, we briefly summarize the solutions to achieve EI model security.

**Defense strategies for poisoning attacks:** Poisoning attacks cause damage to the reliability and usability of models by manipulating and injecting poisoned data. In EI, FL is a commonly used approach for model training, where the adversary usually directly poisons local updates without modifying the training data. Unlike centralized learning frameworks, defenses against poisoning attacks in FL need to take into account that local training data is not visible to the public. One defense idea is to detect malicious clients by distinguishing between toxic and benign updates.

Authors in [42] develop an unsupervised anomaly detection method based on Support Vector Machines (SVM). They introduce a separate validation operation for each potentially malicious local model, to improve anomaly detection accuracy. Differently, authors in [64] propose a weight-based detection scheme. It provides edge nodes with small validation datasets to detect and filter anomalous parameters uploaded by malicious end devices. Based on detection results, edge nodes set appropriate parameter weights to eliminate the effect of pseudo-parameters on the model. In addition, they use DP techniques to provide privacy-preserving measures for sensitive data.

Another approach is to enforce robust aggregation algorithm to counter interference from unknown adversaries. Authors in [35] propose a secure cosine similarity scheme to identify toxic gradients using HE as the underlying technique. The communication overhead is also reduced by replacing the remote communication method with intra-cluster communication. Authors in [36] design a lightweight dual-server secure aggregation protocol that utilizes a third-party server to compute the distance between the model updates from different participants and the average model update. Based on the deviation degrees, the protocol selects the top k nearest participants, and their local updates are considered benign. Authors in [37] use the Euclidean distance of the model to filter malicious updates and subsequently compute the median of the remaining model coordinates to ensure the accuracy of results. In addition, they provide a secure aggregation environment with the help of TEE (i.e., Intel SGX) for global model security. Differently, authors in [65] propose a scoring

TABLE III: Security solutions for learning models.

| Solutions | Ref. | Description | Optimization Metrics |
|---|---|---|---|
| Anomaly detection | [42] | A federal anomaly analysis framework to defend against local model poisoning attacks in distributed ML frameworks. | Security, accuracy, and computation costs |
| | [64] | A FL model combining anomaly detection and DP to resist poisoning attacks and protect privacy. | Security, privacy, and accuracy |
| Robust aggregation | [36] | A novel lightweight privacy-preserving crowdsourced FL scheme to support secure model aggregation. | Security, privacy, and computation costs |
| | [35] | A layered privacy-preserving defense architecture against poisoning attacks in data heterogeneity scenarios. | Security, privacy, accuracy, and communication costs |
| | [37] | An efficient TEE-based aggregation framework for Byzantine robust FL. | Security and accuracy |
| | [65] | A two-stage defense algorithm to analyze local feature patterns of malicious remote updates. | Security and accuracy |
| | [66] | A fast and efficient byzantine robustness algorithm for P2P systems. | Security, efficiency, and accuracy |
| | [67] | A blockchain consensus-based approach to defend against model poisoning attacks. | Security and latency |
| | [44] | A secure and privacy-preserving decentralized learning system. | Security, privacy, and efficiency |
| Adversarial training | [68] | A dynamic defense mechanism to improve EI classification accuracy in adversarial settings. | Security and accuracy |
| GAN | [69] | A decentralized fast vigilance framework for identifying adversarial attacks in IAISs. | Security and latency |

model that employs kernel density estimation to evaluate updates from remote clients. They statistically approximate the optimal threshold to distinguish malicious updates from clean ones.

Instead of relying on a central parameter server, authors in [66] propose a fast and computationally efficient byzantine-robust algorithm for fully decentralized training systems. Their algorithm utilizes a new sequential, memory-assisted, and performance-based criterion for training on logical rings while filtering out byzantine users. Similarly, researchers in [67] and [44] also focus on robust aggregation for fully decentralized training systems. The difference from [66] is that they use blockchain to provide a transparent and secure process.

**Defense strategies for adversarial attacks:** Adversary attacks aim to deceive and mislead the output of a model by making a minor but intentionally designed modification to the input data. Compared to large-scale cloud-based models, compression models running on edge devices typically have fewer parameters and computational resources, and adversarial examples can easily mislead them [68]. Adversarial training enhances the ability to handle attacks by incorporating adversarial examples into the training data. However, these defense measures rely on computationally expensive solutions to generate effective adversarial examples, limiting their applicability in EI.

To achieve lightweight defense, authors in [68] propose a dynamic defense mechanism that combines techniques such as KD, MTD, and Bayesian Stackelberg games to improve

model robustness, especially the classification accuracy of EI in an adversarial setting. Authors in [69] focus on adversarial attacks in Industrial AI Systems (IAISs). The proposed method combines multiple DL models and multiple Conditional Generative Adversarial Networks (CGAN). The CGAN control plane and the DL model data plane operate without the need for complex robustness reinforcement of the original DL model. Author in [79] introduce LanCeX, which is designed to counter compression adversarial attacks in embedded recognition scenarios. LanCeX comprises a detection phase to identify potential adversarial patterns and a data recovery method to mitigate adversarial perturbations in the input data. The proposed defense approach is both universal and lightweight, making it suitable for resource-constrained embedded systems.

*3) Privacy-preserving Solutions for Sensitive Data:* Researchers develop various privacy-preserving schemes aimed at balancing system performance and user privacy protection, providing a solid foundation for reliability and user experience in EI. As shown in Tab. IV, we briefly summarize some solutions.

Authors in [70] propose a lightweight encrypted edge inference systems. They use binary neural networks to reduce resource requirements for deploying models on edge devices, while using secret sharing techniques to provide privacy guarantees. Similarly, authors in [71] use secret sharing based encryption to address privacy-preserving CNN feature extraction for mobile sensing, while significantly reducing the latency and overhead of the end device. Differently, authors in [72]

TABLE IV: Privacy-preserving solutions for sensitive data.

| Solutions | Ref. | Description | Optimization Metrics |
|---|---|---|---|
| Lightweight encryption | [70] | A crypto-neural network inference system at the network edge. | Privacy and latency |
| | [71] | A novel lightweight framework to address privacy issues for mobile sensing. | Privacy, latency, and computation costs |
| | [72] | A lightweight certificate-free multi-user encryption algorithm for mobile devices. | Privacy, communication and computation costs |
| Data desensitization | [73] | A novel hybrid learning approach to enhance privacy protection by extracting special features from sensitive data. | Privacy and latency |
| | [74] | A two-level DP mechanism to enhance privacy protection in cloud-edge-end layered FL frameworks. | Privacy, accuracy, and latency |
| | [75] | A collaborative approach to enhance privacy preservation of inference processes. | Privacy and latency |
| Surrogate data | [76] | A knowledge transfer-based method to trade off data privacy and classification accuracy. | Privacy and accuracy |
| | [77] | A GAN method to generate new data, thereby preserving the privacy of sensitive training data. | Privacy and accuracy |
| | [78] | A novel adversarial sample generation method based on the Firefly algorithm for data protection. | Privacy and latency |

propose a new matching encryption method that provides bilateral access control between the sender and the receiver. The proposed scheme employs a low-power pair-free operation and supports a one-to-many setup to avoid encrypting each receiver's message individually. In addition, they introduce a certificate-less encryption method to address the key escrow problem.

In addition, the risk of privacy leakage is reduced by employing desensitization techniques to process sensitive information, which can be achieved by DP. Authors in [73] propose a privacy-preserving and latency-aware DL framework that addresses privacy and latency issues in EI systems. It introduces inductive learning and a new local DP algorithm that allows edge devices to apply random noise to features extracted from sensitive data before transmitting them to a central server. Authors in [74] utilize a cloud-edge-end hierarchical FL framework to offload the training burden of the device to the nearest edge to improve efficiency, while protecting privacy with a two-level DP mechanism. Differently, authors in [75] establish a collaborative strategy to hide the attributes of the original inputs by distributing CNN feature mappings across multiple heterogeneous IoT devices, thus making it impossible for malicious devices to recover the original data. At the same time, a trade-off between latency and privacy is established.

Privacy protection also can be achieved by replacing sensitive data with similarly distributed proxy data. Authors in [80] propose a defense method, where the model is trained on a different but related dataset, avoiding direct access to original sensitive dataset. Authors in [76] utilize an unprotected model trained on private data and transfer its knowledge to a student model trained on labeled reference data. Authors in [77] employ GANs to generate new data that required for training. To enhance quality of the generated data by GANs, authors

utilize truncation techniques and clustering algorithms during the generation process for different types of data. Similarly, authors in [78] propose a data protection method based on data disturbance and adversarial training. They also introduce a new adversarial sample generation method based on firefly algorithm.

**Lesson 1:** In this subsection, we delve into three aspects of security for edge communication, security for learning models, and privacy protection, providing flexible and balanced security strategies for resource-constrained edge networks. A comprehensive analysis of these solutions reveals several valuable lessons: *i*) In terms of security for edge communications, the use of lightweight encryption and secure communication channels is the key to improve data confidentiality and integrity. Novel authentication methods, such as collaborative authentication, excel in adapting to highly dynamic edge networks and resource limitations; *ii*) Methods such as unsupervised anomaly detection and robust aggregation algorithms are effective in detecting and filtering malicious updates. The use of dynamic defense mechanisms, knowledge distillation, and multi-task learning helps to improve model robustness; *iii*) Methods such as DP, feature extraction and GAN are effective in protecting data privacy while reducing the required computation and communication resources.

However, the trade-off between security and resource efficiency needs to be further optimized, especially given the limitations of edge devices. Additionally, the evolving landscape of cyber threats necessitates adaptive and proactive security strategies. Striking the right balance between security measures and usability remains a persistent challenge in the dynamic context of EI.
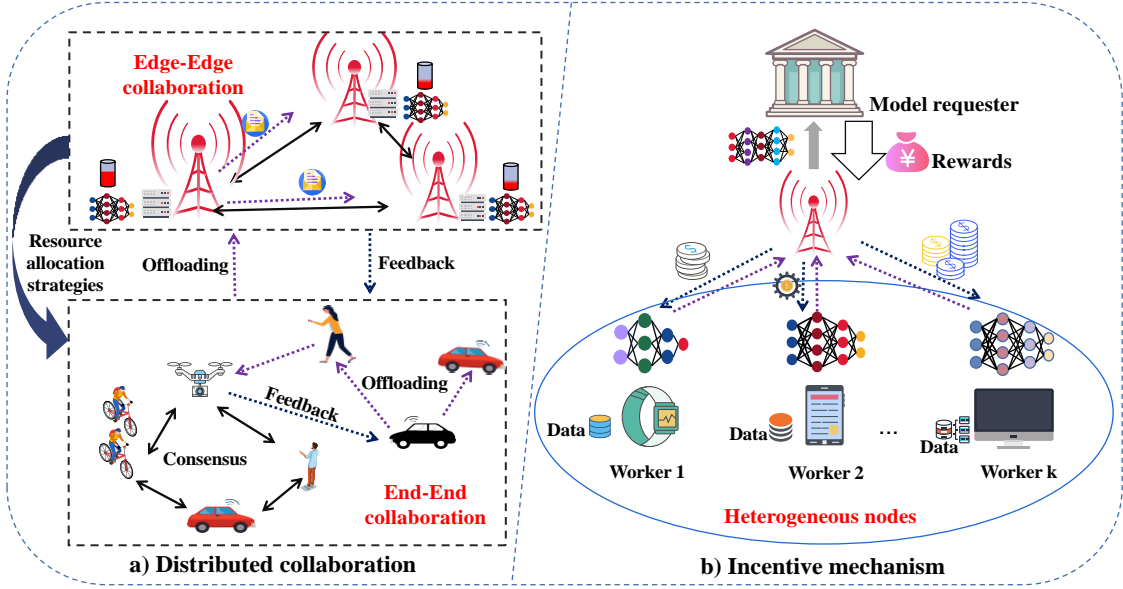
Fig. 5: Illustrative solutions of reliability: a) Distributed collaboration: It allows heterogeneous devices to work together and share knowledge, reducing the risk of single points of failure while increasing resource utilization; and b) Incentive mechanisms: Rewards are given for honest and high-quality resource contributions.

## B. Solutions for Reliability

As discussed in subsection II-C-2, highly dynamic edge environments where the addition of new nodes may change the network connectivity and the departure of nodes may lead to insufficient resources in certain areas of the network. Such topological changes pose challenges to the network stability and reliability. Therefore, the system requires effective measures to adapt to the rapidly changing network environment to ensure user experience, task execution, and system functionality. In this subsection, we detail solutions, including distributed collaboration and incentive mechanisms, to address the issues and thus establish reliability of trustworthy EI, which is shown in Fig. 5.

*1) Distributed Collaboration:* Through distributed collaboration, the EI system is able to build resilient and adaptable characteristics. As shown in Tab.V, we summarize the relevant literature on distributed collaborative EI solutions to achieve reliability, which include: autonomous decision-making, knowledge and model sharing, and decentralized mechanisms.

First, autonomous decision-making assists edge devices in adjusting their behavior to adapt to dynamic environments and task requirements. For instance, authors in [81] focus on resource allocation of multiple UAVs in a ground vehicle network. By using multi-agent RL algorithm, each UAV acts as an intelligent agent with centralized training and distributed execution to collaborate on resource. In addition, an attention mechanism is introduced where agents can further optimize their local models based on information from other agents. Authors in [82] propose an imitation learning based strategy to train distributed agents. Unlike RL, which requires a specific reward function, the imitation learning supports to teach complex tasks to agent with few information, thus turning the training process into fitting an expert demonstration distribution. Authors in [24] propose a multi-level perceptual task offloading framework, in which a vehicle is able to utilize other nearby AVs and Roadside Units (RSUs) to perform collaborative computation for integrated perception of the region of interest. In the proposed approach, the AVs are able to dynamically assign, offload, and execute tasks based on the characteristics and latency requirements of sensing tasks.

In addition, collaborative resource allocation among nodes is a key solution to achieve network reliability and improve QoS. Authors in [83] propose an intelligent resource allocation framework, which trains a DNN to predict resource allocation behavior through self-supervised learning. Subsequently, action prediction is combined with multi-task learning to enhance the performance of DNN. Differently, authors in [84] consider the synergy of computational resources from mobile devices and edge servers to accelerate multi-outlet DNN inference. In the proposed model, bi-directional dynamic programming is used for exit point selection, and then DRL is used to model partitioning and resource allocation strategies.

Second, the overall performance is improved by sharing knowledge, models and experience among collaborative nodes. Authors in [91] present a theoretical framework for real-time evolutionary learning of states in decentralized edge inference network. Specifically, the goal of each agent is to infer a time-varying state in a decentralized manner by using its local observations and messages from other nodes within the communication range. The article tries to provide a communication-efficient coding strategy for generating transmission messages and a sufficient condition for the boundedness of the distributed inference error over time for all agents. In [85], authors introduce a consensus mechanism that utilizes Device-to-Device (D2D) communication to mitigate

TABLE V: Distributed collaboration solutions for reliability.

| Solutions | Ref. | Description | Network Scenarios | Optimization Metrics |
|---|---|---|---|---|
| Autonomous decision-making | [81] | A multi-UAV collaborative system combining multi-agent RL and attention mechanisms for efficient resource allocation. | Multi-UAV assisted vehicular networks | Latency and communication costs |
| | [82] | A lightweight imitation learning-based distributed agent strategy for autonomous cooperation. | Edge networks | Latency |
| | [24] | A perceptual task offloading framework to improve the reliability of autonomous driving. | IoV | Latency |
| | [83] | An intelligent resource allocation framework based on multi-task deep RL algorithm. | IoT | Latency and power consumption |
| | [84] | A novel multi-exit DNN inference acceleration framework. | Edge inference networks | Latency and computation costs |
| Knowledge and model sharing | [85] | A semi-decentralized learning architecture combining device-to-server and D2D communication paradigms. | Edge training networks | Latency, accuracy and energy consumption |
| | [86] | A FL framework for edge server collaboration based on decentralized consensus. | Edge training network | Latency and accuracy |
| | [87] | A fragmented distributed learning approach for IoT devices. | IoT | Latency, computation costs, and energy consumption |
| Decentralized collaboration | [88] | A decentralized DRL-based resource allocation and task scheduling approach. | UAV-MEC networks | Latency and energy consumption |
| | [89] | A collaborative DRL approach for resource allocation that ensures long-term energy savings. | MEC networks | Energy efficiency |
| | [90] | An efficient decentralized GNN-based training algorithm. | Decentralized communication networks | Communication costs |

model disagreement and improve resource efficiency. This work can be considered as a learning approach in the middle star topology between traditional FL and fully decentralized architectures, resulting in a new semi-decentralized learning architecture. Similarly, authors in [86] devise a joint optimization method to facilitate collaborative learning among a large number of edge devices in a wide region. By utilizing multiple aggregators to mitigate concerns about single point of failure, the approach is more scalable than the traditional FL framework.

Different from [85] and [86], authors in [87] propose fragmented learning in resource-constrained IoT edge devices. The method splits the training operation for each data point into atomic operations and executes them on Fog Nodes (FNs). An iterative implementation of fragmentation learning enables FNs to transfer partially learned weights to the next appropriate FN for further training, and repeats the process until the training is complete. By analyzing main parameters and using a greedy heuristic algorithm, the algorithm selects optimal FNs for the training operation, which tries to reduce the probability of interruptions caused by device failures.

Finally, decentralized mechanisms for decision making and task assignment ensure that nodes can work together without centralized control. Authors in [88] propose a multi-agent DRL approach for decentralized implementation, where multiple UAVs collaborate to determine their computational and communication strategies. Similarly, authors in [89] present a fully decentralized multi-agent DRL algorithm for autonomous re-

source allocation in heterogeneous mobile edge networks. The algorithm has a multi-actor shared criticism architecture and a regional training distributed execution framework, which aims to stabilize model training and reduce information exchange.

Graph Neural Networks (GNNs) have attracted much academic attention for its ability to efficiently process graph data, adapt to the dynamic nature of wireless networks, and enable decentralized management and control. Authors in [90] propose a personalized training algorithm based on graph attention to address the problems caused by non-independently identically and distributed data in traditional distributed learning. The algorithm allows each agent to train a local model that personalizes data by learning the specific weights of different neighboring nodes without prior knowledge of the graph structure or the data distribution of neighboring nodes. However, the fading of wireless channels makes GNN affected by the information exchange among neighbors [100]. Therefore, authors in [100] try to enhance the robustness of decentralized GNNs in the inference phase through two new re-transmission mechanisms.

*2) Incentive Mechanisms:* Appropriate incentive mechanisms can inspire edge nodes to actively participate in collaboration and contribute their own resources, which helps to mitigate the impact of node withdrawal, thus improving the overall reliability and adaptability of the system. In the following, we introduce incentive mechanisms based on node contributions and blockchain, which are summarized in Tab. VI.

TABLE VI: Incentivization solutions for reliability.

| Categories | Ref. | Solutions | | | | | Description |
|---|---|---|---|---|---|---|---|
| | | Auction Theory | Game Theory | Contract Theory | Smart Contract | DRL | |
| Incentivization via node contributions | [92] | $\checkmark$ | $\times$ | $\times$ | $\times$ | $\times$ | A framework for integrating quality estimation, reverse auction incentives, and automatic weighted aggregation for high-quality model training. |
| | [93] | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ | $\times$ | A two-tier incentive mechanism that considers local data contribution and energy consumption of training nodes. |
| | [94] | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ | $\times$ | A two-tier resource allocation and incentive mechanism for decentralized learning-based systems. |
| | [95] | $\times$ | $\times$ | $\times$ | $\times$ | $\checkmark$ | A DRL-based automatic pricing strategy to incentivize nodes to contribute computational resources. |
| | [96] | $\times$ | $\checkmark$ | $\times$ | $\times$ | $\times$ | A Bayesian game-based incentive mechanism to encourage nodes to contribute their data and computational resources in FL. |
| Incentivization via blockchain | [97] | $\times$ | $\checkmark$ | $\times$ | $\checkmark$ | $\times$ | A dynamic incentive model combining evolutionary game theory and smart contracts to encourage users to participate in data sharing. |
| | [98] | $\times$ | $\times$ | $\checkmark$ | $\times$ | $\times$ | A contract-based incentive mechanism to encourage nodes to participate in training and share resources, and store the reputation of nodes via blockchain. |
| | [45] | $\times$ | $\checkmark$ | $\times$ | $\checkmark$ | $\times$ | A blockchain-based FL incentive mechanism to balance system overhead and model performance. |
| | [99] | $\times$ | $\checkmark$ | $\times$ | $\times$ | $\times$ | A two-stage Stackelberg game approach to optimize resource allocation in blockchain-based FL systems. |

("$\checkmark$" if the research satisfies the solution, "$\times$" if not)

**Incentivization via node contributions:** The quality of model updates can vary greatly due to a variety of factors such as training data volume, data quality, and data distribution. As a result, some studies utilize them to measure customer contributions. Authors in [92] employ a reverse auction model to incentivize high-quality and low-cost computing nodes to participate in training process.

A single-layer incentive mechanism may be not efficient in motivating all parties. To address this challenge, researchers in [93] and [94] propose a two-layer incentive mechanism. In the lower layer, the size of reward pool and data allocation of training nodes are determined based on data volume, data quality, and privacy budget provided by each data owner. In the upper layer, profits are allocated to training nodes based on their marginal contributions to the model publisher, i.e., the impact on the global model performance.

Instead of focusing solely on a single dimension, such as node data quality and quantity, which may not guarantee long term system efficiency and stability, some studies address situations where nodes possess multi-dimensional attributes, such as heterogeneity communication resources and computational capabilities [45]. Authors in [95] discuss incentives for using DRL to balance the trade-off between latency and payment in edge learning. At the beginning of each training round, the aggregator publishes the price of each Edge Node (EN) and the EN determine the computational resources they provide for model training. The incentives are presented in the form

of a hierarchical game, using an actor-critic network model to overcome the challenges of the dynamic environment and the access to parameter information of ENs. Similarly, authors in [96] consider the balance among rewards, costs, edge communication and computing capabilities, allowing nodes to determine their contributions of local data and computational resources in each training round.

**Incentivization via blockchain:** Blockchain-based incentives create a trusted and transparent environment for collaborative EI. It is able to record and preserve the historical reputation of nodes, preventing interference from unreliable participants [101]. Authors in [97] combine reputation-based and payment-based incentives, using "credibility coins" as encrypted cryptocurrency for data transactions. At the same time, they introduce a dynamic incentive model based on evolutionary game theory to analyze user interactions and the stability of strategies. Authors in [98] integrate reputation and contract theory to ensure fair rewards. They utilize blockchain for secure reputation management of training nodes, providing non-repudiation and anti-tampering properties in a decentralized manner. These reputation-based approaches assess the trustworthiness and contributions of participating nodes, either based on their impacts on model performance or through social interactions, to ensure integrity and reliability of learning.

Moreover, authors in [45] leverage blockchain-based incentive mechanisms to optimize the balance between system overhead and model performance. They achieve this by

compensating relevant nodes with different resources, while also enhancing data privacy through mechanisms embedded in smart contracts. This approach establishes a credible, faster, and transparent resource trading system. Similarly, in [99], the authors formulate the resource allocation problem within blockchain-based learning as a two-stage Stackelberg game. They assist model owners in reward allocation and clients in determining their computational resources for both model training and task mining.

**Lesson 2:** With solutions such as distributed collaboration and incentives, it is possible to build an intelligent system at the network edge with a high degree of flexibility and reliability to adapt to dynamic communication environments and collaboration among heterogeneous nodes. This lays the foundation for achieving sustainable QoS. Distributed collaboration enables that nodes to communicate with each other, share knowledge, and tolerate single point of failure. In particular, the combination of blockchain and smart contracts helps to ensure the credibility and transparency of the incentives.

However, when implementing and optimizing distributed collaboration and incentives, it is necessary to ensure fairness among nodes and guard against free-riding behavior. This includes avoiding over-centralization of rewards while ensuring that nodes share resources honestly. In addition, the manageability of distributed and decentralized networks is challenging compared to centralized architectures and requires a combination of technical, institutional and social considerations.

### C. Solutions for Transparency

The transparency of decision-making processes allows users to trust the model, while allowing system administrators to monitor and adjust the model behavior. According to the description in subsection III-C-3, although deep models used for edge training and inference have undergone compression and processing, it is still a complex black-box model. Therefore, achieving intrinsic interpretability of models becomes extremely difficult. In contrast to model intrinsic interpretability, model decision interpretability focuses on explaining specific decisions without attention to the interpretability of the entire model structure and parameters. In subsection III-B, we introduce several widely used interpretation techniques, which can be broadly classified into three categories: vision-based interpretation, feature-based interpretation, and alternative model-based interpretation. These approaches focus on revealing key factors that lead to particular outputs of the model, thereby improving the understanding and interpretation of model decisions. As shown in Tab. VII, we briefly summarize the literature related to interpretable solutions of decision-making.

**Visual-based explanations** offer an intuitive and visual representation of decision-making processes and underlying patterns of models. For instance, authors in [102] refine the CAM technique to highlight relevant portions of EEG signals associated with mental states in the vehicle drowsiness monitoring system. Authors in [25] utilize the guided-Grad-CAM method to provide real-time explanations with high-resolution activation maps for multi-modal DL models. Different from

CAM, authors in [53] utilize attention modules for fine-grained spatial localization, generating accurate heatmaps. Furthermore, research in [103] combines micro and macro interpretation modules to explain the failure cases of object detection models in autonomous driving systems. By extracting and visualizing features of CNNs and providing spatio-temporal information, this method aims to assist model developers in understanding, fine-tuning, and developing models.

**Feature attribution** methods provide explanations by analyzing the importance of features that affect the output. LIME and SHAP are two commonly used feature-based post-hoc interpretation methods. Authors in [51] use DeepSHAP method to provide local and global explanations for DL-based intrusion detection systems in IoV networks. Differently, authors in [47] combine LIME and SHAP to explain the importance of clinical features and genotype in warfarin daily dose prediction models. This combination also performs global and local interpretations to help healthcare practitioners understand and trust model predictions. In the global interpretation, the ranked importance and SHAP interpreter produce a ranking of feature importance across the entire dataset. In local interpretation, LIME and SHAP interpreters show the effect of features on the output of models run on specific samples.

**Surrogate models** approximate the behavior of black box models and provide explanations. Authors in [104] propose a tree regularization technique to approximate the complex decision boundaries of deep models. By optimizing the surrogate models, domain experts can gain a good understanding of black box models' behavior. In contrast, authors in [105] propose a fuzzy rule-based local agent model that is capable of providing a model-independent interpretation. In addition, the fuzzy rule-based model allows one to strike a balance between prediction accuracy and the number of rules. Once established, the granularity of information in the fuzzy rule backparts can be interpreted, transforming lengthy sequences of numbers into concise and detailed descriptors, and endowing the generated fuzzy rules with a high degree of interpretability.

**Lesson 3:** In this subsection, we summarize interpretable implementations of lightweight models deployed at the network edge to resolve the complex black-box processing problem. This not only enhances transparency in EI systems, but also ensures that the decisions are rational and ethically responsible. In terms of decision interpretability, vision-based interpretation may be more intuitive for certain tasks, while feature-based interpretation is more suitable for scenarios that require detailed characterization. Multiple interpretation methods can also be combined with each other to provide users with a comprehensive interpretation of the model. It is important to note that ensuring the security of interpretation techniques is also an important direction for future research. Attackers can utilize interpretable techniques to detect vulnerabilities, especially in risk-sensitive scenarios [106]. In addition, there is a lack of scientific evaluation systems to assess interpretation methods.

### D. Solution for Sustainability

In this subsection, we delve into solutions for the trade-off between energy consumption and QoS, and the collection and

TABLE VII: Transparency solutions based on decision interpretability.

| Solutions | Ref. | Description | Technologies |
|---|---|---|---|
| Visual-based explanations | [102] | A specific interpretation technique designed for EEG signal detection models to reveal localized regions of input signals that are important for prediction. | CAM |
| | [25] | A novel interpretable multi-modal DL model to recognize important biomarkers in brain images. | Grad-CAM |
| | [53] | An interpretable lightweight CNN design based on an attention mechanism. | Attention mechanism |
| | [103] | A visual attention model generates fine-grained causal visual attention heat maps to explain the steering control of AVs. | Attention mechanism |
| Feature-based explanations | [51] | A DL framework based on the SHAP approach to provide local and global interpretation for transparent and robust intrusion detection. | SHAP |
| | [47] | An interpretable framework for integrated predictive modeling based on LIME and SHAP to explain the importance of relevant features. | SHAP, LIME |
| Surrogate model-based explanations | [104] | A tree regularization technique for approximating complex decision boundaries of DL models, to provide interpretability and maintain predictive accuracy. | Regularization |
| | [105] | A fuzzy local alternative model to improve the interpretability of learning model results. | Fuzzy rule |

storage of high-quality data. These solutions not only address the performance imbalance caused by edge nodes with limited resources, but also help ensure their sustainable operations.

*1) Trade-offs Between Energy Consumption and QoS:* Communication and computation delays directly affect the responsiveness of AI applications and is particularly important for tasks that require real-time decision making and fast feedback, such as AVs and industrial automation. By considering both latency and energy consumption, systems can minimize resource consumption while delivering high-quality services, thus achieving sustainability in'edge environments. As shown in Tab. VIII, we provide a brief overview of solutions on sustainability.

Authors in [107] consider both proximity constraints and capacity constraints of edge servers, aiming to maximize the number of users served while minimizing system energy consumption. They propose an energy saving policy that is dynamically formulated over time, deciding the state in which the edge server is woken up or put into hibernation. A heterogeneous framework is proposed in [26], which organizes various heterogeneous resources and is designed to achieve efficient utilization and intelligent scheduling, as well as provide a reliable and sustainable environment. The node scheduler is responsible for recording the status and information of all edge nodes, including the current workload, hardware specifications, and software settings. When an AI task request is received, the system find the most suitable edge node for execution to minimize the processing energy consumption while meeting the latency requirements.

In many scenarios, fixed edge servers are often overloaded in hotspots and expensive to deploy in remote areas, making the above approaches difficult to realize in specific contexts. Therefore, UAV-based EI architecture becomes a viable solution, especially for areas where service resources are scarce. For example, flexible UAVs equipped with AIGC servers enable users to access AIGC services with ultra-low latency and high reliability [28]. Therefore, the energy-efficient design of UAVs is important to ensure their sustainable operation. Authors in [108] focus on sustainable services in UAV-assisted MEC networks. This article aims to minimize energy consumption of UAVs by jointly optimizing flight trajectories, resource allocation, and task scheduling with fairness. Differently, authors in [109] not only consider energy consumption of UAVs, but also focus on uplink transmission of mobile users. They use DRL algorithms to solve the joint optimization problem of UAV motion control, mobile user association and power control.

In deploying large model services to the network edge, the above resource scheduling and offloading strategies can be utilized to balance the overall service latency and device energy consumption. Furthermore, some studies focus on joint optimization of latency, accuracy and energy consumption of model training and inference processes. For instance, authors in [110] present a multi-agent RL collaborative inference scheme that enables each device to select the best DNN division point and collaborative edge based on the number of images, channel conditions and previous inference performance. The Q-values of neighboring devices are used to accelerate the optimization of the inference strategy through learning experience exchange, thus reducing energy consumption while ensuring inference latency and quality. Differently, authors in [111] propose auto-encoder based intermediate feature compression to reduce the communication overhead in collaborative inference thereby achieving fast and energy-efficient edge inference.

Unlike [110] and [111], studies in [112] and [113] focus on transformer models in the edge inference process. These studies are important for efficient and accurate reasoning based on transformer models, such as mobile AIGC. Authors in [112] decompose large vision transformers into multiple smaller

TABLE VIII: Balanced energy consumption and QoS solutions for sustainability.

| Ref. | Description | Network scenarios | Optimization Metrics |
|---|---|---|---|
| [107] | A dynamic scheduling strategy for energy-efficient MEC services. | MEC networks | Energy consumption and user coverage |
| [26] | An online scheduling strategy for energy-efficient AI services. | IoT | Energy consumption and latency |
| [108] | An efficient resource allocation and fair task offloading scheme for sustainable services in UAV-assisted MEC networks. | UAV-MEC networks | Energy consumption and latency |
| [109] | A real-time resource allocation scheme based on DRL algorithm to reduce energy consumption and system latency in UAV networks. | UAV-MEC networks | Latency, energy consumption and service coverage |
| [110] | An energy efficient edge collaborative AI model inference framework based on multi-agent RL. | Edge inference networks | Energy consumption, latency, and accuracy |
| [111] | A lightweight feature compression method for fast and energy-efficient collaborative inference. | Edge inference networks | Latency and energy consumption |
| [112] | A collaborative inference framework for transformer model decomposition based on knowledge distillation. | Edge inference networks | Energy consumption, latency, and accuracy |
| [113] | A comprehensive framework for implementing edge deployment based on the transformer model. | Edge inference networks | Energy consumption, latency, and accuracy |
| [114] | A flexible weight quantification approach for energy-efficient FL. | Edge training networks | Energy consumption, latency, and accuracy |
| [115] | An energy-efficient FL based on CPU-GPU heterogeneous computing. | Edge training networks | Energy consumption, latency, and accuracy |
| [116] | A novel framework to minimize energy consumption of FL in MEC networks, where the availability of UEs is uncertain. | Edge training networks | Energy consumption and accuracy |
| [117] | A dynamic device scheduling algorithm with energy-aware FL frameworks. | Edge training networks | Energy consumption and accuracy |

models to be deployed at edge devices to reduce inference latency and energy consumption. In addition, the loss of model accuracy caused by decomposition is reduced by an algorithm based on knowledge distillation. Authors in [113] introduce multiple frameworks to enable efficient deployment of transformer architectures to EI platforms. Specifically, ProTran and FlexiBERT 2.0 are used for modeling accuracy as well as latency, energy consumption, and hardware measurements. EdgeTran employs alternative models derived from ProTran and FlexiBERT 2.0 to obtain the best-performing model-device pairs and performs the optimization of the output model through block-level growth and castration techniques.

Moreover, as for collaborative training frameworks across edge devices such as FL, local computation and frequent communication may overwhelm energy-constrained mobile devices. A possible solution is to develop multiple model compression methods such as pruning and quantization. For example, in the study of [114], an iterative algorithm for jointly determining weight quantization and spectrum resource allocation strategies among devices is proposed.

Some studies address the balance between learning performance and energy consumption by resources management. Authors in [115] achieve energy-efficient FL in wireless networks by considering the allocation of both computational and communication resources on four dimensions: bandwidth

allocation, time partitioning, CPU-GPU workload partitioning, and CPU-GPU frequency scaling. Authors in [116] leverage complex interactions among environmental contextual information (e.g., workload, the amount of available computational resources, and data quality) to select available User Equipments (UEs) and propose an approximation algorithm to find the suitable aggregator location. Similarly, authors in [117] also introduce an energy-aware device scheduling algorithm. Unlike [116], this article employs an analog gradient aggregation scheme, which aims to achieve more efficient device scheduling by aggregating local updates over the same time-frequency resource blocks.

*2) High-quality Data Collection and Storage:* High-quality data can fuel the training, inference, and optimization of AI models. Therefore, data collection and sustainable data storage are important to ensure reliable network services.

**Data collection:** Mobile Crowd Sensing (MCS) utilizes smart devices and sensors that are widely distributed at the network edge to cellularize data to support various intelligent applications. However, it becomes difficult to directly assess the quality of collected data without any prior knowledge. To address the issue, a data quality-based MCS incentive mechanism is designed in [118], which aims to enable long-term effective contributions by paying participants who provide high-quality data. Authors in [27] propose a cost and quality-

aware data collection scheme for edge-assisted vehicle crowd-sensing systems. They employ adaptive clustering and online sensing parameter tuning to avoid unnecessary data collection while ensuring reliable and timely data uploads.

Authors in [119] utilize blockchain for data quality control. The blockchain records key information about the employment relationship, participants and consensus nodes to ensure that the data is undeniable and tamper-proof. Accurate reward payments are realized based on user data quality, truth results, reward rules and employment relationships.

In addition, multi-modal data fusion helps to integrate information from different modalities, including vision, speech, and sensor data. Authors in [120] combine modal data fusion techniques and GAN to build a cross-modal data generator. The generator can generate long-term time series data from spatial-temporal modal data, and then replace missing values with the generated data.

**Data storage:** A reliable and scalable storage solution ensures data persistence and availability. Authors in [121] focus on ensuring cache fairness among peers in edge environments. They propose approximation algorithms and distributed algorithms to determine cache nodes. Further, fairness metrics are utilized to achieve continuous caching decisions over time. Authors in [122] explore various criteria, such as data prevalence and proximity to edge processing functions, to strategically allocate diverse classes of raw and processed data at the network edge. They provide a lightweight mechanism to identify data types, data generation sources, and supported processing functions.

Due to the limited storage capacity of edge servers, reducing data redundancy is of great importance to improve the storage utilization of edge servers [123]. Deduplication methods are widely used to reduce data redundancy in storage systems. Authors in [123] use integer programming and Lagrangian relaxation methods, as well as an improved subgradient approach to solve the balanced deduplication problem. In the work presented in [124], the simultaneous consideration of file popularity, file similarity, and server reliability aims to enhance the availability of popular files while minimizing unnecessary space redundancy. The article underscores the adverse effects of server unreliability on storage hit rates and addresses this issue by implementing similarity-aware hierarchical clustering algorithms.

**Lessons 4:** In this subsection, we focus on two key perspectives, efficient energy utilization and data quality, to facilitate sustainable and long-term service quality in resource-constrained edge environments. This is crucial for the efficient deployment and inference of AI models at the network edge. By recognizing that optimizing one factor may lead to sacrifices in other aspects, striking a balance between energy consumption and service quality empowers the system to provide high-quality services while minimizing resource consumption, thereby enhancing the overall sustainability of the EI system.

Additionally, due to the varying quality of data, including noise and biases, there can be impacts on the accuracy of large models. Therefore, multi-modal data fusion and GANs can provide rich and realistic data representations for various tasks.

This approach is also helpful in learning complex relationships among data to support advanced applications. Nevertheless, research on these technologies is still in its infancy.

## V. Research Challenges and Open Issues

Research on trustworthy EI is still in its infancy. In order to achieve the vision of secure, reliable, transparent, and sustainable trustworthy EI, there are many issues and research directions that need to be further explored. Inspired by existing solutions, in this section, we discuss some research challenges and open issues in the field of trustworthy EI.

### A. Security Enhancement

*1) Security and Privacy in D2D-based Decentralized Paradigm:* Decentralized EI relies solely on communication and consensus among end devices to process tasks, thus overcoming the dependence on central nodes while reducing the high cost of communications with edge servers. However, in this scenario, there are several security and privacy risks, such as: *i)* Decentralized entity control: Different nodes in a decentralized system may be controlled by distinct administrative entities, leading to potential inconsistencies in the security level of nodes. A vulnerable node may become a target of attackers, posing a threat to the overall system security; *ii)* Widespread semi-honest nodes: They can introduce instability into the system, and attempt to steal data or send malicious data during collaborative learning to disrupt system availability; *iii)* Topology diversity: The topology of a distributed system can be highly complex, with diverse relationships among nodes, complicating security analysis and management. Therefore, robust decentralized consensus algorithms need to be developed as well as optimized topology management strategies for distributed systems, while constrained resources of devices need to be considered.

*2) Deployment of Zero-Trust Security at Network Edge:* Incorporating the zero-trust principle in EI architectures is an effective way to address both known and unknown threats. However, due to limited edge resources, it faces the following challenges and open issues: *i)* Computational and storage resource requirements: zero-trust model requires continuous monitoring and verification, imposing an excessive burden on the limited computational and storage resources of edge devices, and a lightweight implementation needs to be found; *ii)* Low-latency requirements: EI services typically require low latency, but multiple validation steps introduced by the zero-trust model increase latency; *iii)* Diversity of edge devices: Considering the diversity of devices in edge environments, including different hardware and software specifications, the realization of zero-trust needs to be deployed consistently and efficiently on various types of devices.

*3) Generative Model-based Security:* In the field of EI, generative models such as GANs and Variational AutoEncoder (VAE) are widely used to improve the security of systems. These generative models are capable of capturing unusual patterns or behaviors and detecting inconsistencies with normal behavior by learning normal data distributions, so that potential security threats can be detected in a timely

manner. In adversarial defense of images, these methods have achieved remarkable results. However, there are challenges in migrating these approaches to communication network defense for EI. Communications in EI systems involve a wider range of data types and complex communication patterns, thus how to effectively apply generative models to detect anomalous behaviors is an issue that requires in-depth research. On the other hand, by using intelligent game-adversarial techniques, it is possible to realize the induced detection computation of EI models by "generative robots". By automating the generation of large test sets, these "generative robots" are able to detect model weaknesses and security issues. This approach provides a cutting-edge research direction for discovering and providing feedback on security issues in EI systems.

### B. Autonomous Collaboration for Edge Co-inference

Autonomous collaboration among edge nodes is essential for dynamic edge co-inference, but its implementation poses challenges in distributed training and execution. In distributed control scenarios, edge nodes can make independent decisions and execute tasks in a dynamic environment while collaborating to achieve common goals. To realize such autonomous decision-making and distributed collaboration, current research typically uses distributed RL.

Fully distributed RL is particularly challenging, because it needs to consider not only the interaction between individual agents and the environment, but also the interplay among multiple agents. Correspondingly, various challenges are introduced, including: *i*) Changes in agent strategies: It can lead to environmental instability, since the behavior of one agent affects that of others; *ii*) Distributed training and reward feedback: Distributed training requires individual agents to receive separate reward feedback. Decomposing feedback from the environment into rewards for each agent and quantifying contributions of each agent to teamwork are rather complex; *iii*) Curse of dimensionality: When the number of agents increases, the learning process faces challenges such as the curse of dimensionality, resulting in a significant increase in computational complexity.

### C. Large Models in Trustworthy EI

Large models, especially LLMs, are deployed at the network edge as model collaboration control center of EI system, which not only coordinates multiple intelligent applications, but also ensures their consistency, accuracy, and real-time performance in handling complex tasks. This deployment helps to build trustworthy EI system and provide users with secure and reliable intelligent services. Specifically, EI systems will integrate multiple intelligent applications, such as AI assistants that are not only limited to smart home management, but also include real-time environment monitoring, user communication, and so on. These complex tasks are beyond the solution capability of a single AI model. Smart speakers, home control screens, smart phones, and other IoT devices are expected to become the interaction portal of Jarvis style smart housekeeper. Therefore, considering that immediacy, reliability, privacy, and computational capability, it is crucial to deploy large models as a
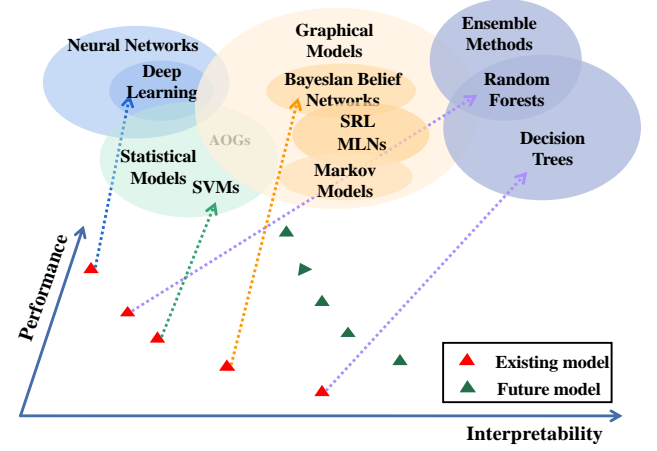


Fig. 6: Accuracy of interpretable models [126].

collaborative control center for models on edge devices. However, this work is still in its infancy, and further breakthroughs are needed in compression techniques to achieve efficient deployment and support tens of billions of parametric model inference at the network edge. Key challenges include low-latency inference that matches endpoint computing capability, the impact of accuracy loss on user QoS, and the sustainability challenge of energy consumption.

### D. Trade-off Between Interpretability and Accuracy at The Network Edge

In EI, balancing accuracy and interpretability is an important issue [125]. Some models can achieve high accuracy, but it is difficult to explain their internal mechanisms, limiting their applications in various fields. Other models with high interpretability may sacrifice accuracy, which also affects their effectiveness in practical applications. Fig. 6 shows the relationship between accuracy and interpretability of ML models.

This trade-off issue involves several challenging factors: *i*) "Interpretability" is difficult to be defined and measured because different application fields and scenarios have distinct requirements; *ii*) In EI application fields like healthcare, protecting data privacy is crucial, limiting the use of models with strong interpretability; *iii*) Some models with strong interpretability may be too simplistic to handle large-scale and high-dimensional data, potentially sacrificing accuracy; *iv*) Some models may require many parameters or complex structures to achieve high accuracy, which are hard for interpretation; *v*) In scenarios involving high-dimensional and complex data, even models with strong interpretability may struggle to explain their internal mechanisms and decision-making processes, necessitating the development of novel interpretative methods. In conclusion, achieving a balance between accuracy and interpretability in trustworthy EI is a complex task, requiring the considerations of multiple factors and practical application requirements.

## E. Data Governance Issues in Trustworthy EI

EI integrates sensing, communication, and computation, which corresponds to the collection, transmission, and utilization of data. High-quality data resources are crucial for models to accurately understand and respond to various scenarios, especially generative AI large models that require a large amount of high-quality data to enhance content accuracy. However, multi-modal data contains rich societal knowledge. For example, digital data, as a potentially huge resource, is difficult to apply for big model training due to its simple expression form and lack of linguistic features. In addition, utilizing generated data to train AI models will lead to model degradation. Specifically, the generated data obtained by sampling the output of the generative model will lose some information. This situation can gradually accumulate and ultimately lead to a distribution of the generated data that bears little resemblance to the real data. As a result, models trained with such generated data will lead to degraded model performance and unreliable inference results [127]. Therefore, the efficient fusion of multi-modal data and high-quality synthetic data will become the difficulty of breakthrough.

Furthermore, ensuring visibility into data use and sharing can help to improve users' understanding of how data are used and increase their sense of control and trust in data privacy. However, achieving such visibility requires strong technical support, including real-time monitoring, logging and permission tracking. Overcoming these technical challenges requires significant investment in resources and technology development. Finding the suitable balance between improving visibility while balancing the user's right to privacy and the protection of sensitive information, such as trade secrets, is a complex task.

## VI. Conclusion

We summarize and discuss the development, solutions and challenges in a large number of related literatures on trustworthy EI. First, we discuss the definition, characteristics, and architecture of trustworthy EI. We also sort out four important issues in the implementation of trustworthy EI and summarize the key enabled techniques to achieve the trustworthiness of EI. Subsequently, we conduct a comprehensive investigation from different aspects, i.e., balanced security and privacy protection, reliability, transparency, and sustainability of trustworthy EI. Finally, we discuss the relevant research challenges and open issues toward achieving trustworthy EI. We hope this survey provides an effective guideline that can inspire researchers to advance trustworthy EI.

## References

[1] H. Djigal, J. Xu, L. Liu, and Y. Zhang, "Machine and deep learning for resource allocation in multi-access edge computing: A survey," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 4, pp. 2449–2494, 2022.

[2] E. Baccour, N. Mhaisen, A. A. Abdellatif, A. Erbad, A. Mohamed, M. Hamdi, and M. Guizani, "Pervasive AI for IoT applications: A survey on resource-efficient distributed artificial intelligence," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 4, pp. 2366–2418, 2022.

[3] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.

[4] Z. Zhang, B. Yang, B. Wang, and B. Li, "GrowSP: Unsupervised semantic segmentation of 3D point clouds," in *Proc. IEEE/CVF CVPR*, 2023, pp. 17 619–17 629.

[5] O. M. Manyar, Z. McNulty, S. Nikolaidis, and S. K. Gupta, "Inverse reinforcement learning framework for transferring task sequencing policies from humans to robots in manufacturing applications," in *Proc. IEEE ICRA*, 2023, pp. 849–856.

[6] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.

[7] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1738–1762, 2019.

[8] Y. Zhao, Z. Yang, X. He, X. Cai, X. Miao, and Q. Ma, "Trine: Cloud-edge-device cooperated real-time video analysis for household applications," *IEEE Transactions on Mobile Computing*, vol. 22, no. 8, pp. 4973–4985, 2023.

[9] S. Baker and W. Xiang, "Artificial intelligence of things for smarter healthcare: A survey of advancements, challenges, and opportunities," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 2, pp. 1261–1293, 2023.

[10] S. V. Balkus, H. Wang, B. D. Cornet, C. Mahabal, H. Ngo, and H. Fang, "A survey of collaborative machine learning using 5G vehicular communications," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 2, pp. 1280–1303, 2022.

[11] Z. Ning, H. Hu, X. Wang, L. Guo, S. Guo, G. Wang, and X. Gao, "Mobile edge computing and machine learning in the Internet of unmanned aerial vehicles: A survey," *ACM Computing Surveys, doi:10.1145/3604933*, 2023.

[12] K. B. Letaief, Y. Shi, J. Lu, and J. Lu, "Edge artificial intelligence for 6G: Vision, enabling technologies, and applications," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 1, pp. 5–36, 2022.

[13] D. Kaur, S. Uslu, K. J. Rittichier, and A. Durresi, "Trustworthy artificial intelligence: A review," *ACM Computing Surveys*, vol. 55, no. 2, pp. 1–38, 2022.

[14] H. Liu, Y. Wang, W. Fan, X. Liu, Y. Li, S. Jain, Y. Liu, A. Jain, and J. Tang, "Trustworthy AI: A computational perspective," *ACM Transactions on Intelligent Systems and Technology*, vol. 14, no. 1, pp. 1–59, 2022.

[15] B. Li, P. Qi, B. Liu, S. Di, J. Liu, J. Pei, J. Yi, and B. Zhou, "Trustworthy AI: From principles to practices," *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–46, 2023.

[16] M. M. H. Shuvo, S. K. Islam, J. Cheng, and B. I. Morshed, "Efficient acceleration of deep learning inference on resource-constrained edge devices: A review," *Proceedings of the IEEE*, vol. 111, no. 1, pp. 42–91, 2023.

[17] M. S. Murshed, C. Murphy, D. Hou, N. Khan, G. Ananthanarayanan, and F. Hussain, "Machine learning at the network edge: A survey," *ACM Computing Surveys*, vol. 54, no. 8, pp. 1–37, 2021.

[18] Y. Shi, K. Yang, T. Jiang, J. Zhang, and K. B. Letaief, "Communication-efficient Edge AI: Algorithms and systems," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 4, pp. 2167–2191, 2020.

[19] M. A. Al-Garadi, A. Mohamed, A. K. Al-Ali, X. Du, I. Ali, and M. Guizani, "A survey of machine and deep learning methods for Internet of things (IoT) security," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 1646–1685, 2020.

[20] B. Mao, J. Liu, Y. Wu, and N. Kato, "Security and privacy on 6G network edge: A survey," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 2, pp. 1095–1127, 2023.

[21] H. Du, R. Zhang, D. Niyato, J. Kang, Z. Xiong, D. I. Kim, X. S. Shen, and H. V. Poor, "Exploring collaborative distributed diffusion-based AI-generated content (AIGC) in wireless networks," *IEEE Network, doi:10.1109/MNET.006.2300223*, 2023.

[22] H. Zhao, H. Chen, F. Yang, N. Liu, H. Deng, H. Cai, S. Wang, D. Yin, and M. Du, "Explainability for large language models: A survey," *arXiv preprint*, 2023. [Online]. Available: https://doi.org/10.48550/arXiv.2309.01029

[23] A. Kumar, V. Capraro, and M. Perc, "The evolution of trust and trustworthiness," *Journal of the Royal Society Interface*, vol. 17, no. 169, p. 20200491, 2020.

[24] Z. Xiao, J. Shu, H. Jiang, G. Min, H. Chen, and Z. Han, "Perception task offloading with collaborative computation for autonomous driving," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 2, pp. 457–473, 2022.

[25] W. Hu, X. Meng, Y. Bai, A. Zhang, G. Qu, B. Cai, G. Zhang, T. W. Wilson, J. M. Stephen, V. D. Calhoun, and Y.-P. Wang, "Interpretable multimodal fusion networks reveal mechanisms of brain cognition," *IEEE Transactions on Medical Imaging*, vol. 40, no. 5, pp. 1474–1483, 2021.

[26] S. Zhu, K. Ota, and M. Dong, "Green AI for IIoT: Energy efficient intelligent edge computing for industrial Internet of things," *IEEE Transactions on Green Communications and Networking*, vol. 6, no. 1, pp. 79–88, 2021.

[27] L. Liu, L. Wang, Z. Lu, Y. Liu, W. Jing, and X. Wen, "Cost-and-quality aware data collection for edge-assisted vehicular crowdsensing," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 5, pp. 5371–5386, 2022.

[28] M. Xu, H. Du, D. Niyato, J. Kang, Z. Xiong, S. Mao, Z. Han, A. Jamalipour, D. I. Kim, V. Leung *et al.*, "Unleashing the power of edge-cloud generative AI in mobile networks: A survey of AIGC services," *arXiv preprint*, 2023. [Online]. Available: https://doi.org/10.48550/arXiv.2303.16129

[29] Y. Lin, Z. Gao, H. Du, D. Niyato, J. Kang, R. Deng, and X. S. Shen, "A unified blockchain-semantic framework for wireless edge intelligence enabled Web 3.0," *IEEE Wireless Communications, doi:10.1109/MWC.018.2200568*, 2023.

[30] M. Xu, W. C. Ng, W. Y. B. Lim, J. Kang, Z. Xiong, D. Niyato, Q. Yang, X. Shen, and C. Miao, "A full dive into realizing the edge-enabled metaverse: Visions, enabling technologies, and challenges," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 656–700, 2023.

[31] J.-H. Cho, D. P. Sharma, H. Alavizadeh, S. Yoon, N. Ben-Asher, T. J. Moore, D. S. Kim, H. Lim, and F. F. Nelson, "Toward proactive, adaptive defense: A survey on moving target defense," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 1, pp. 709–745, 2020.

[32] Y. Zhou, G. Cheng, and S. Yu, "An SDN-enabled proactive defense framework for DDoS mitigation in IoT networks," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 5366–5380, 2021.

[33] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. Vincent Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.

[34] A. Acar, H. Aksu, A. S. Uluagac, and M. Conti, "A survey on homomorphic encryption schemes: Theory and implementation," *ACM Computing Surveys*, vol. 51, no. 4, pp. 1–35, 2018.

[35] X. Chen, H. Yu, X. Jia, and X. Yu, "Apfed: Anti-poisoning attacks in privacy-preserving heterogeneous federated learning," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 5749–5761, 2023.

[36] Z. Zhang, L. Wu, C. Ma, J. Li, J. Wang, Q. Wang, and S. Yu, "LSFL: A lightweight and secure federated learning scheme for edge computing," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 365–379, 2022.

[37] L. Zhao, J. Jiang, B. Feng, Q. Wang, C. Shen, and Q. Li, "Sear: Secure and efficient aggregation for byzantine-robust federated learning," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 5, pp. 3329–3342, 2021.

[38] P. Fu, J. Wu, X. Lin, and A. Shen, "ZTEI: Zero-trust and edge intelligence empowered continuous authentication for satellite networks," in *Proc. IEEE GLOBECOM*, 2022, pp. 2376–2381.

[39] T. Wang, N. Huang, Y. Wu, J. Gao, and T. Q. S. Quek, "Latency-oriented secure wireless federated learning: A channel-sharing approach with artificial jamming," *IEEE Internet of Things Journal*, vol. 10, no. 11, pp. 9675–9689, 2023.

[40] Z. Xu and M. Baykal-Gürsoy, "Power allocation for cooperative jamming against a strategic eavesdropper over parallel channels," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 846–858, 2023.

[41] H. Yang, Z. Xiong, J. Zhao, D. Niyato, L. Xiao, and Q. Wu, "Deep reinforcement learning-based intelligent reflecting surface for secure wireless communications," *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, pp. 375–388, 2021.

[42] S. Shi, C. Hu, D. Wang, Y. Zhu, and Z. Han, "Federated anomaly analytics for local model poisoning attack," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 2, pp. 596–610, 2021.

[43] X. Wang, Z. Ning, L. Guo, S. Guo, X. Gao, and G. Wang, "Mean-field learning for edge computing in mobile blockchain networks," *IEEE Transactions on Mobile Computing*, vol. 22, no. 10, pp. 5978–5994, 2023.

[44] M. Xu, Z. Zou, Y. Cheng, Q. Hu, D. Yu, and X. Cheng, "SPDL: A blockchain-enabled secure and privacy-preserving decentralized learning system," *IEEE Transactions on Computers*, vol. 72, no. 2, pp. 548–558, 2023.

[45] X. Wang, Y. Zhao, C. Qiu, Z. Liu, J. Nie, and V. C. Leung, "InFEdge: A blockchain-based incentive mechanism in hierarchical federated learning for end-edge-cloud communications," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 12, pp. 3325–3342, 2022.

[46] M. T. Ribeiro, S. Singh, and C. Guestrin, ""Why should I trust you?" Explaining the predictions of any classifier," in *Proc. ACM SIGKDD*, 2016, pp. 1135–1144.

[47] Y. Zhang, C. Xie, L. Xue, Y. Tao, G. Yue, and B. Jiang, "A post-hoc interpretable ensemble model to feature effect analysis in warfarin dose prediction for chinese patients," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 2, pp. 840–851, 2022.

[48] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. NeurIPS*, vol. 30, 2017, pp. 4768–4777.

[49] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, "From local explanations to global understanding with explainable AI for trees," *Nature Machine Intelligence*, vol. 2, no. 1, pp. 56–67, 2020.

[50] J. Cheng, M. Gao, J. Liu, H. Yue, H. Kuang, J. Liu, and J. Wang, "Multimodal disentangled variational autoencoder with game theoretic interpretability for glioma grading," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 2, pp. 673–684, 2022.

[51] A. Oseni, N. Moustafa, G. Creech, N. Sohrabi, A. Strelzoff, Z. Tari, and I. Linkov, "An explainable deep learning framework for resilient intrusion detection in IoT-enabled transportation networks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 1, pp. 1000–1014, 2023.

[52] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE/CVF ICCV*, 2017, pp. 618–626.

[53] Y. S. Jeon, K. Yoshino, S. Hagiwara, A. Watanabe, S. T. Quek, H. Yoshioka, and M. Feng, "Interpretable and lightweight 3-D deep learning model for automated ACL diagnosis," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 7, pp. 2388–2397, 2021.

[54] J. Zhang, F. Zhang, X. Huang, and X. Liu, "Leakage-resilient authenticated key exchange for edge artificial intelligence," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 6, pp. 2835–2847, 2021.

[55] B. Gong, C. Guo, C. Guo, C. Guo, Y. Sun, M. Waqas, and S. Chen, "SLIM: A secure and lightweight multi-authority attribute-based signcryption scheme for IoT," *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 1299–1312, 2024.

[56] R. Kaewpuang, M. Xu, W. Y. B. Lim, D. Niyato, H. Yu, J. Kang, and X. S. Shen, "Cooperative resource management in quantum key distribution (QKD) networks for semantic communication," *IEEE Internet of Things Journal, doi:10.1109/JIOT.2023.3301033*, 2023.

[57] B. Li, T. Shi, W. Zhao, and N. Wang, "Reinforcement learning-based intelligent reflecting surface assisted communications against smart attackers," *IEEE Transactions on Communications*, vol. 70, no. 7, pp. 4771–4779, 2022.

[58] Q. He, S. Fang, T. Wang, Y. Liu, S. Zhao, and Z. Lu, "Proactive anti-eavesdropping with trap deployment in wireless networks," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 1, pp. 637–649, 2023.

[59] T. Zhang, Y. Huo, Q. Gao, L. Ma, Y. Wu, and R. Li, "Cooperative physical layer authentication with reputation-inspired collaborator selection," *IEEE Internet of Things Journal*, vol. 10, no. 24, pp. 22 165–22 181, 2023.

[60] H. Fang, Z. Xiao, X. Wang, L. Xu, and L. Hanzo, "Collaborative authentication for 6G networks: An edge intelligence based autonomous approach," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 2091–2103, 2023.

[61] J. Cui, N. Liu, Q. Zhang, D. He, C. Gu, and H. Zhong, "Efficient and anonymous cross-domain authentication for IIoT based on blockchain," *IEEE Transactions on Network Science and Engineering*, vol. 10, no. 2, pp. 899–910, 2023.

[62] F. Tong, X. Chen, K. Wang, and Y. Zhang, "CCAP: A complete cross-domain authentication based on blockchain for Internet of things," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 3789–3800, 2022.

[63] X. Wang, H. Zhu, Z. Ning, L. Guo, and Y. Zhang, "Blockchain intelligence for Internet of vehicles: Challenges and solutions," *IEEE Communications Surveys & Tutorials, doi:10.1109/COMST.2023.3305312*, 2023.

[64] J. Zhou, N. Wu, Y. Wang, S. Gu, Z. Cao, X. Dong, and K.-K. R. Choo, "A differentially private federated learning model against poisoning attacks in edge computing," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 3, pp. 1941–1958, 2023.

[65] X. Li, Z. Qu, S. Zhao, B. Tang, Z. Lu, and Y. Liu, "LoMar: A local defense against poisoning attack on federated learning," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 1, pp. 437–450, 2023.

[66] A. R. Elkordy, S. Prakash, and S. Avestimehr, "Basil: A fast and Byzantine-resilient approach for decentralized training," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2694–2716, 2022.

[67] D. C. Nguyen, S. Hosseinalipour, D. J. Love, P. N. Pathirana, and C. G. Brinton, "Latency optimization for blockchain-empowered federated learning in multi-server edge computing," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 12, pp. 3373–3390, 2022.

[68] Y. Qian, Y. Guo, Q. Shao, J. Wang, B. Wang, Z. Gu, X. Ling, and C. Wu, "EI-MTD: moving target defense for edge intelligence against adversarial attacks," *ACM Transactions on Privacy and Security*, vol. 25, no. 3, pp. 1–24, 2022.

[69] G. Li, K. Ota, M. Dong, J. Wu, and J. Li, "DeSVig: Decentralized swift vigilance against adversarial attacks in industrial artificial intelligence systems," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 5, pp. 3267–3277, 2020.

[70] X. Liu, B. Wu, X. Yuan, and X. Yi, "Leia: A lightweight cryptographic neural network inference system at the edge," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 237–252, 2022.

[71] K. Huang, X. Liu, S. Fu, D. Guo, and M. Xu, "A lightweight privacy-preserving CNN feature extraction framework for mobile sensing," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 3, pp. 1441–1455, 2021.

[72] N. Yang, C. Tang, and D. He, "A lightweight certificateless multi-user matchmaking encryption for mobile devices: Enhancing security and performance," *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 251–264, 2024.

[73] J. Tchaye-Kondi, Y. Zhai, J. Shen, and L. Zhu, "Privacy-preserving offloading in edge intelligence systems with inductive learning and local differential privacy," *IEEE Transactions on Network and Service Management*, vol. 20, no. 4, pp. 5026–5037, 2023.

[74] Y. Guo, F. Liu, T. Zhou, Z. Cai, and N. Xiao, "Privacy vs. efficiency: Achieving both through adaptive hierarchical federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 4, pp. 1331–1342, 2023.

[75] E. Baccour, A. Erbad, A. Mohamed, M. Hamdi, and M. Guizani, "RL-DistPrivacy: Privacy-aware distributed deep inference for low latency IoT systems," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 4, pp. 2066–2083, 2022.

[76] V. Shejwalkar and A. Houmansadr, "Membership privacy for machine learning models through knowledge transfer," in *Proc. AAAI*, vol. 35, no. 11, 2021, pp. 9549–9557.

[77] L. Hu, J. Li, G. Lin, S. Peng, Z. Zhang, Y. Zhang, and C. Dong, "Defending against membership inference attacks with high utility by GAN," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 3, pp. 2144–2157, 2023.

[78] P. Zhang, Y. Wang, N. Kumar, C. Jiang, and G. Shi, "A security- and privacy-preserving approach based on data disturbance for collaborative edge computing in social IoT systems," *IEEE Transactions on Computational Social Systems*, vol. 9, no. 1, pp. 97–108, 2022.

[79] Z. Xu, F. Yu, C. Liu, and X. Chen, "LanCeX: A versatile and lightweight defense method against condensed adversarial attacks in image and audio recognition," *ACM Transactions on Embedded Computing Systems*, vol. 22, no. 1, pp. 1–24, 2022.

[80] H. Huang, W. Luo, G. Zeng, J. Weng, Y. Zhang, and A. Yang, "DAMIA: leveraging domain adaptation as a defense against membership inference attacks," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 5, pp. 3183–3199, 2021.

[81] Y. Wang, Y. He, F. R. Yu, Q. Lin, and V. C. M. Leung, "Efficient resource allocation in multi-UAV assisted vehicular networks with security constraint and attention mechanism," *IEEE Transactions on Wireless Communications*, vol. 22, no. 7, pp. 4802–4813, 2023.

[82] Z. Ning, H. Chen, E. C. H. Ngai, X. Wang, L. Guo, and J. Liu, "Lightweight imitation learning for real-time cooperative service migration," *IEEE Transactions on Mobile Computing*, vol. 23, no. 2, pp. 1503–1520, 2024.

[83] J. Chen, S. Chen, Q. Wang, B. Cao, G. Feng, and J. Hu, "IRAF: A deep reinforcement learning approach for collaborative mobile edge computing IoT networks," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 7011–7024, 2019.

[84] F. Dong, H. Wang, D. Shen, Z. Huang, Q. He, J. Zhang, L. Wen, and T. Zhang, "Multi-exit DNN inference acceleration based on multi-dimensional optimization for edge intelligence," *IEEE Transactions on Mobile Computing*, vol. 22, no. 9, pp. 5389–5405, 2023.

[85] F. P.-C. Lin, S. Hosseinalipour, S. S. Azam, C. G. Brinton, and N. Michelusi, "Semi-decentralized federated learning with cooperative D2D local model aggregations," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 12, pp. 3851–3869, 2021.

[86] Z. Zhang, Z. Gao, Y. Guo, and Y. Gong, "Scalable and low-latency federated learning with cooperative mobile edge networking," *IEEE Transactions on Mobile Computing*, vol. 23, no. 1, pp. 812–822, 2024.

[87] P. K. Deb, A. Mukherjee, D. Singh, and S. Misra, "Loop-the-loops: Fragmented learning over networks for constrained IoT devices," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 1, pp. 316–327, 2022.

[88] S. Hwang, H. Lee, J. Park, and I. Lee, "Decentralized computation offloading with cooperative UAVs: Multi-agent deep reinforcement learning perspective," *IEEE Wireless Communications*, vol. 29, no. 4, pp. 24–31, 2022.

[89] Y. Xiao, Y. Song, and J. Liu, "Collaborative multi-agent deep reinforcement learning for energy-efficient resource allocation in heterogeneous mobile edge computing networks," *IEEE Transactions on Wireless Communications, doi:10.1109/TWC.2023.3335597*, 2023.

[90] Z. Tian, Z. Zhang, Z. Yang, R. Jin, and H. Dai, "Distributed learning over networks with graph-attention-based personalization," *IEEE Transactions on Signal Processing*, vol. 71, pp. 2071–2086, 2023.

[91] Z. Liu, A. Conti, S. K. Mitter, and M. Z. Win, "Communication-efficient distributed learning over networks–part I: Sufficient conditions for accuracy," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 4, pp. 1081–1101, 2023.

[92] Y. Deng, F. Lyu, J. Ren, Y.-C. Chen, P. Yang, Y. Zhou, and Y. Zhang, "Improving federated learning with quality-aware user incentive and auto-weighted model aggregation," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 12, pp. 4515–4529, 2022.

[93] J. S. Ng, W. Y. B. Lim, Z. Xiong, X. Cao, D. Niyato, C. Leung, and D. I. Kim, "A hierarchical incentive design toward motivating participation in coded federated learning," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 1, pp. 359–375, 2021.

[94] W. Y. B. Lim, J. S. Ng, Z. Xiong, J. Jin, Y. Zhang, D. Niyato, C. Leung, and C. Miao, "Decentralized edge intelligence: A dynamic resource allocation framework for hierarchical federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 3, pp. 536–550, 2022.

[95] Y. Zhan and J. Zhang, "An incentive mechanism design for efficient edge learning by deep reinforcement learning approach," in *Proc. IEEE INFOCOM*, 2020, pp. 2489–2498.

[96] M. Hu, W. Yang, Z. Luo, X. Liu, Y. Zhou, X. Chen, and D. Wu, "AutoFL: A bayesian game approach for autonomous client participation in federated edge learning," *IEEE Transactions on Mobile Computing*, vol. 23, no. 1, pp. 194–208, 2024.

[97] Y. Chen, Y. Zhang, S. Wang, F. Wang, Y. Li, Y. Jiang, L. Chen, and B. Guo, "DIM-DS: Dynamic incentive model for data sharing in federated learning based on smart contracts and evolutionary game theory," *IEEE Internet of Things Journal*, vol. 9, no. 23, pp. 24 572–24 584, 2022.

[98] J. Kang, Z. Xiong, D. Niyato, S. Xie, and J. Zhang, "Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 10 700–10 714, 2019.

[99] Z. Wang, Q. Hu, R. Li, M. Xu, and Z. Xiong, "Incentive mechanism design for joint resource allocation in blockchain-based federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 5, pp. 1536–1547, 2023.

[100] M. Lee, G. Yu, and H. Dai, "Decentralized inference with graph neural networks in wireless communication systems," *IEEE Transactions on Mobile Computing*, vol. 22, no. 5, pp. 2582–2598, 2023.

[101] Z. Ning, S. Sun, X. Wang, L. Guo, G. Wang, X. Gao, and R. Y. Kwok, "Intelligent resource allocation in mobile blockchain for privacy and security transactions: A deep reinforcement learning based approach," *Science China Information Sciences*, vol. 64, no. 6, p. 162303, 2021.

[102] J. Cui, Z. Lan, O. Sourina, and W. Müller-Wittig, "EEG-based cross-subject driver drowsiness recognition with an interpretable convolutional neural network," *IEEE Transactions on Neural Networks and Learning Systems, doi:10.1109/TNNLS.2022.3147208*, 2022.

[103] J. Wang, Y. Li, Z. Zhou, C. Wang, Y. Hou, L. Zhang, X. Xue, M. Kamp, X. Zhang, and S. Chen, "When, where and how does it fail? A spatial-temporal visual analytics approach for interpretable object detection in autonomous driving," *IEEE Transactions on Visualization and Computer Graphics, doi:10.1109/TVCG.2022.3201101*, 2022.

[104] M. Wu, M. Hughes, S. Parbhoo, M. Zazzi, V. Roth, and F. Doshi-Velez, "Beyond sparsity: Tree regularization of deep models for interpretability," in *Proc. AAAI*, vol. 32, no. 1, 2018.

[105] X. Zhu, D. Wang, W. Pedrycz, and Z. Li, "Fuzzy rule-based local surrogate models for black-box model explanation," *IEEE Transactions on Fuzzy Systems*, vol. 31, no. 6, pp. 2056–2064, 2023.

[106] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju, "Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods," in *Proc. AIES*, 2020, pp. 180–186.

[107] G. Cui, Q. He, X. Xia, F. Chen, and Y. Yang, "Eesaver: Saving energy dynamically for green multi-access edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 7, pp. 2155–2166, 2023.

[108] M. Zhao, W. Li, L. Bao, J. Luo, Z. He, and D. Liu, "Fairness-aware task scheduling and resource allocation in UAV-enabled mobile edge computing networks," *IEEE Transactions on Green Communications and Networking*, vol. 5, no. 4, pp. 2174–2187, 2021.

[109] J. Chen, X. Cao, P. Yang, M. Xiao, S. Ren, Z. Zhao, and D. O. Wu, "Deep reinforcement learning based resource allocation in multi-UAV-aided MEC networks," *IEEE Transactions on Communications*, vol. 71, no. 1, pp. 296–309, 2023.

[110] Y. Xiao, L. Xiao, K. Wan, H. Yang, Y. Zhang, Y. Wu, and Y. Zhang, "Reinforcement learning based energy-efficient collaborative inference for mobile edge computing," *IEEE Transactions on Communications*, vol. 71, no. 2, pp. 864–876, 2023.

[111] Z. Hao, G. Xu, Y. Luo, H. Hu, J. An, and S. Mao, "Multi-agent collaborative inference via DNN decoupling: Intermediate feature compression and edge learning," *IEEE Transactions on Mobile Computing*, vol. 22, no. 10, pp. 6041–6055, 2023.

[112] G. Xu, Z. Hao, Y. Luo, H. Hu, J. An, and S. Mao, "DeViT: Decomposing vision transformers for collaborative inference in edge devices," *IEEE Transactions on Mobile Computing, doi:10.1109/TMC.2023.3315138*, 2023.

[113] S. Tuli and N. K. Jha, "EdgeTran: Device-aware co-search of transformers for efficient inference on mobile edge platforms," *IEEE Transactions on Mobile Computing, doi:10.1109/TMC.2023.3328287*, 2023.

[114] R. Chen, L. Li, K. Xue, C. Zhang, M. Pan, and Y. Fang, "Energy efficient federated learning over heterogeneous mobile devices via joint design of weight quantization and wireless transmission," *IEEE Transactions on Mobile Computing*, vol. 22, no. 12, pp. 7451–7465, 2023.

[115] Q. Zeng, Y. Du, K. Huang, and K. K. Leung, "Energy-efficient resource management for federated edge learning with CPU-GPU heterogeneous computing," *IEEE Transactions on Wireless Communications*, vol. 20, no. 12, pp. 7947–7962, 2021.

[116] Z. Xu, D. Li, W. Liang, W. Xu, Q. Xia, P. Zhou, O. F. Rana, and H. Li, "Energy or accuracy? Near-optimal user selection and aggregator placement for federated learning in MEC," *IEEE Transactions on Mobile Computing, doi:10.1109/TMC.2023.3262829*, 2023.

[117] Y. Sun, S. Zhou, Z. Niu, and D. Gündüz, "Dynamic scheduling for over-the-air federated edge learning with energy constraints," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 1, pp. 227–242, 2022.

[118] D. Peng, F. Wu, and G. Chen, "Data quality guided incentive mechanism design for crowdsensing," *IEEE Transactions on Mobile Computing*, vol. 17, no. 2, pp. 307–319, 2018.

[119] J. An, Z. Wang, X. He, X. Gui, J. Cheng, and R. Gui, "PPQC: A blockchain-based privacy-preserving quality control mechanism in crowdsensing applications," *IEEE/ACM Transactions on Networking*, vol. 30, no. 3, pp. 1352–1367, 2022.

[120] M. Kang, R. Zhu, D. Chen, X. Liu, and W. Yu, "CM-GAN: A cross-modal generative adversarial network for imputing completely missing data in digital industry," *IEEE Transactions on Neural Networks and Learning Systems, doi:10.1109/TNNLS.2023.3284666*, pp. 1–10, 2023.

[121] Y. Huang, X. Song, F. Ye, Y. Yang, and X. Li, "Fair and efficient caching algorithms and strategies for peer data sharing in pervasive edge computing environments," *IEEE Transactions on Mobile Computing*, vol. 19, no. 4, pp. 852–864, 2020.

[122] A.-C. Nicolaescu, S. Mastorakis, and I. Psaras, "Store edge networked data (SEND): A data and performance driven edge storage framework," in *Proc. IEEE INFOCOM, doi:10.1109/INFOCOM42981.2021.9488804*, 2021.

[123] R. Luo, H. Jin, Q. He, S. Wu, and X. Xia, "Enabling balanced data deduplication in mobile edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 5, pp. 1420–1431, 2023.

[124] J. Xia, G. Cheng, L. Luo, D. Guo, P. Lv, and B. Sun, "The doctrine of mean: Realizing deduplication storage at unreliable edge," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 10, pp. 2811–2826, 2023.

[125] A.-M. Nussberger, L. Luo, L. E. Celis, and M. J. Crockett, "Public attitudes value interpretability but prioritize accuracy in artificial intelligence," *Nature Communications*, vol. 13, no. 1, p. 5821, 2022.

[126] D. Gunning, M. Stefik, J. Choi, T. Miller, S. Stumpf, and G.-Z. Yang, "XAI-Explainable artificial intelligence," *Science Robotics*, vol. 4, no. 37, p. eaay7120, 2019.

[127] I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, and R. Anderson, "The curse of recursion: Training on generated data makes models forget," *arXiv preprint*, 2023. [Online]. Available: https://doi.org/10.48550/arXiv.2305.17493