

# Secure and Trustworthy Artificial Intelligence-Extended Reality (AI-XR) for Metaverses

Adnan Qayyum<sup>1,2</sup>, Muhammad Atif Butt<sup>2</sup>, Hassan Ali<sup>2</sup>, Muhammad Usman<sup>3</sup>, Osama Halabi<sup>4</sup>, Ala Al-Fuqaha<sup>5</sup>,  
Qammer H. Abbasi<sup>1</sup>, Muhammad Ali Imran<sup>1</sup> and Junaid Qadir<sup>4\*</sup>

<sup>1</sup> James Watt School of Engineering, University of Glasgow, Glasgow, United Kingdom

<sup>2</sup> Information Technology University (ITU), Punjab, Lahore, Pakistan

<sup>3</sup> Glasgow Caledonian University, Glasgow, United Kingdom

<sup>4</sup> Qatar University, Doha, Qatar

<sup>5</sup> Hamad Bin Khalifa University, Doha, Qatar

**Abstract**—Metaverse is expected to emerge as a new paradigm for the next-generation Internet, providing fully immersive and personalised experiences to socialize, work, and play in self-sustaining and hyper-spatio-temporal virtual world(s). The advancements in different technologies like augmented reality, virtual reality, extended reality (XR), artificial intelligence (AI), and 5G/6G communication will be the key enablers behind the realization of AI-XR metaverse applications. While AI itself has many potential applications in the aforementioned technologies (e.g., avatar generation, network optimization, etc.), ensuring the security of AI in critical applications like AI-XR metaverse applications is profoundly crucial to avoid undesirable actions that could undermine users' privacy and safety, consequently putting their lives in danger. To this end, we attempt to analyze the security, privacy, and trustworthiness aspects associated with the use of various AI techniques in AI-XR metaverse applications. Specifically, we discuss numerous such challenges and present a taxonomy of potential solutions that could be leveraged to develop secure, private, robust, and trustworthy AI-XR applications. To highlight the real implications of AI-associated adversarial threats, we designed a metaverse-specific case study and analyzed it through the adversarial lens. Finally, we elaborate upon various open issues that require further research interest from the community.

## I. INTRODUCTION

In the recent era, metaverse technology is rapidly emerging and there are a lot of potential applications that can benefit from these developments such as healthcare, industry, business, etc. While there is no single agreed-upon definition of a metaverse [1], the metaverse is a convergence of physical, augmented, and virtual reality and provides a powerfully immersive experience to users by allowing them to seamlessly interact with the real and virtual (computer-simulated) world. The term “metaverse” is a combination of two terms, i.e., “meta” which means transcending, and “universe” which refers to the physical universe and the current virtual world. This is the basic definition of the term metaverse, nevertheless, the literature shows that it does not have a unified definition [1].

Metaverse allows the creation of shared virtual space by connecting all virtual worlds through the Internet, where digital avatars (i.e., users) can communicate and interact with each other similar to the physical world. Key backbone technologies in the metaverse include artificial intelligence (AI) and extended reality (XR) that leverage different technological developments such as virtual reality (VR), augmented reality (AR), and mixed reality (MR). In addition, similar to the current Internet, metaverse will leverage other concomitant technologies like information and communication technologies (ICT), 5G, and 6G, but metaverse will provide a qualitatively different experience to its users by enabling real-life-like 3-D experiences through the incorporation of aforementioned technologies.

Metaverse allows for the digitalization of traditional brick-and-mortar institutions and businesses—it will be possible to develop virtual markets, digital lands, and digital infrastructure, which can be bought and sold using blockchain and non-fungible tokens (NFTs), which are non-interchangeable units of data stored on a blockchain. Metaverse can be a game changer in terms of the impact of its potential applications due to the greater immersion, involvement, and personalization possible due to AI-XR. This is the prime reason various corporations have shown great interest in the idealization of the metaverse and are making big bets on developing their own AI-XR-based metaverse ecosystems.

Metaverse is receiving increasing traction from numerous major tech companies worldwide such as Facebook (which is recently rebranded with the name “Meta”), Microsoft, Google, and Amazon. In addition, the widespread adoption of the metaverse is evident in the infusion of billions of dollars of investments by these companies in an attempt to achieve great technological transformation. However, despite such huge traction of the metaverse and its potential to transform existing ecosystems like healthcare, there are numerous challenges associated with the use of AI in the metaverse that may hinder their seamless adoption by end users in the longer term. In addition, in the backdrop of recent technologically induced social issues, there is a palpable lack of trust and confidence

\*Corresponding author: Junaid Qadir (jqadir@qu.edu.qa)

TABLE I

COMPARISON OF OUR PAPER WITH EXISTING SURVEYS AND REVIEW PAPERS THAT ARE FOCUSED ON ANALYZING PRIVACY AND SECURITY OF AI-XR METAVERSE APPLICATIONS. (LEGEND: S → SECURITY; P → PRIVACY; R → REGULATORY; T → TRUSTWORTHY; √ → COVERED; × → NOT COVERED; ≈ → PARTIALLY COVERED)

Year	Authors	Focused Area	General Issues			ML Related Issues & Solutions				Background & Applications	Open Issues
			S	P	R	S	P	T	XAI		
2018	Falchuk et al. [2]	Privacy issues and solutions for digital footprints in metaverse games.	√	√	×	×	×	×	×	×	×
2020	Guzman et al. [3]	Analyzed privacy and security in MR from data-centric perspective.	√	√	≈	×	×	×	×	×	√
2021	Ning et al. [4]	General focused survey on metaverse with partial discussion on privacy and security issues.	≈	≈	×	×	×	×	×	×	√
2021	Pietro et al. [5]	Discussed general privacy and security issues in metaverse applications.	√	√	×	×	×	×	×	×	√
2022	Huynh et al. [6]	Discussed potential applications of AI in various metaverse applications.	≈	≈	×	×	×	×	×	×	≈
2022	Zhao et al. [7]	Security & privacy issues and solutions for four dimensions: communication, user information, scenario, and goods.	√	√	×	×	×	×	×	×	×
2022	Wang et al. [8]	Presented general (non ML-associated) security and privacy related challenges for different metaverse applications.	√	√	√	×	×	×	×	×	√
	This Paper	ML-associated security, privacy, and trustworthiness challenges and solutions for AI-XR metaverse applications.	√	√	√	√	√	√	√	√	√

in such technologies.

Since technology can be used both ways (i.e., for good and harm), it is vital that governments, corporations, and society at large seriously consider ethical and moral issues. There are many ethical questions about privacy, security, transparency, accountability, democracy, freedom of speech, and anonymity that technology alone cannot answer. Some specific concerns related to how AI-XR-based metaverse applications will impact humanity are: (1) how AI-XR-based metaverse applications will impact and promote human values and human rights? how will AI-XR-based metaverse promote social welfare and not cause harm to society at large; (3) how can we regulate critical applications of metaverse like healthcare? (4) How do we align the commercial and technical imperatives of AI-XR metaverse applications with human values and promote a moral vision and character development? (5) How do we ensure that AI-XR metaverse application developers do not exploit or manipulate their users? Keeping in mind the aforementioned questions, in this paper, we present an analysis of the security, privacy, and trust issues associated with the use of AI-XR in metaverse applications.

*Contributions of this Paper:* To the best of our knowledge this paper is the first attempt towards analyzing the challenges associated with the use of different AI techniques in AI-XR metaverse applications. The comparison of this paper with existing survey and review articles that are focused on analyzing security and privacy aspects of AI-XR metaverse applications is presented in Table I. In the summary, the following are the salient contributions of this paper.

- 1) We highlight various issues that arise with the use of AI in metaverse applications that mainly include security, privacy, and trustworthiness.
- 2) We present a taxonomy of different potential solutions that can be used to realize secure, robust, safe, and trustworthy AI-XR applications.
- 3) We identify various ML-based use cases across different layers of metaverse architecture and highlight several ML-associated vulnerabilities at each layer.
- 4) We present a case study to highlight the real threat of AI-based vulnerabilities by considering a prospective metaverse application design scenario.
- 5) We elaborate upon various open issues that require further development.

*Organization of this Paper:* The rest of the paper is orga-

nized as follows. Section II presents relevant background. The discussion of challenges related to security, safety, privacy, and trust is presented in Section III. The taxonomy of different vulnerabilities associated with the use of ML in AI-XR metaverse applications is discussed in Section IV. Different potential solutions that can ensure security, privacy, safety, and trust in ML applications are discussed in Section V. Various open issues that require further research attention are highlighted in Section VI. Finally, we conclude the paper in Section VII. The organization of the paper is depicted in Figure 1.

## II. BACKGROUND

### A. Metaverse: An Introduction

Before understanding the concept of the metaverse, it is very important to understand the related concepts that are described below.

- *Virtual Reality (VR):* In VR, the users achieve an immersive experience by donning a VR headset that allows them to enter into a virtual (computer-simulated) world thus completely blanking out the real world. The key objective of enabling immersion in VR is to provide high fidelity user interaction to give him the feeling that the virtual world is real [9]. Prime examples of VR include Facebook Oculus and HTC VIVE VR headrests. VR has a wide range of applications but a VR headset is required to enter the digital world.
- *Augmented Reality (AR):* In AR, the users obtain an immersive experience by blending the virtual (digital) and the real world and projecting digital content (text, images, and sounds) onto the real world. Unlike VR, AR can be realized without special equipment (like a headset) through the use of smartphones, implants, glasses, or contact lenses that are used to overlay digital content on top of the real world.
- *Mixed Reality (MR):* MR is a hybrid term that is used to refer to the conjunction of virtual and real worlds to produce new environments and experiences, where physical and digital objects co-exist and interact in real-time (it is an enhanced form of AR). Microsoft HoloLens headrest is an example MR headrest.
- *Extended Reality (XR):* XR is an umbrella term that encompasses and subsumes VR, AR, and MR. It covers all the future realities that might emerge from these

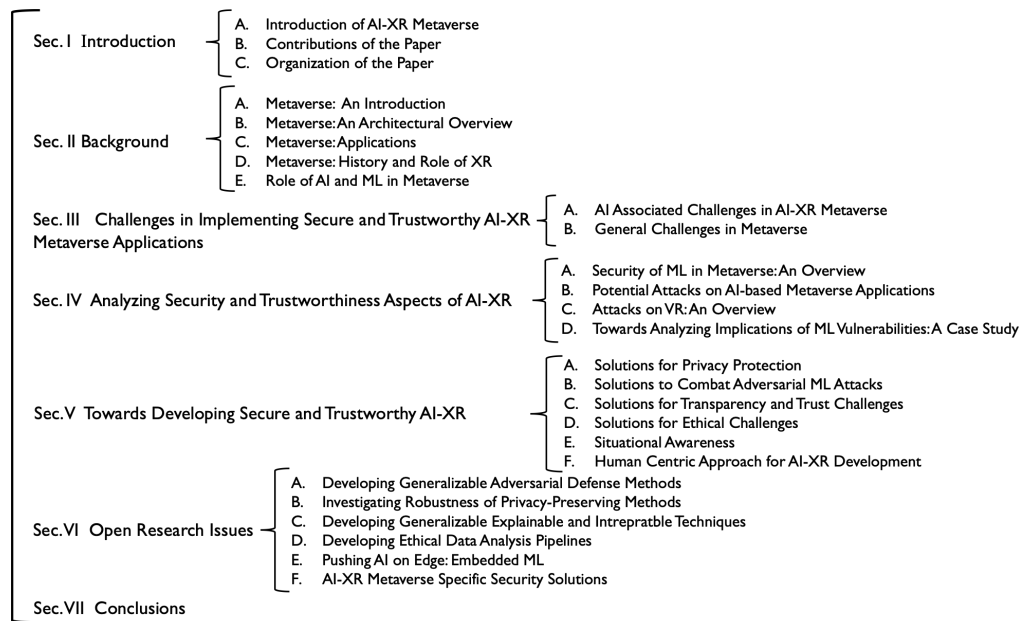


Fig. 1. Organization of the paper.

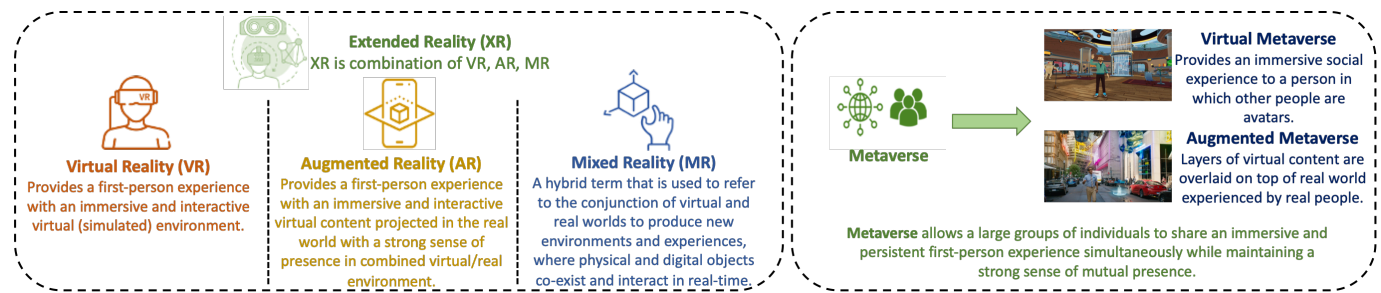


Fig. 2. An overview of different concepts related to metaverse that include VR, AR, MR, XR, virtual metaverse, and augmented metaverse.

technologies. XR is predicted to become a \$209 billion market by 2022.

Immersive first-person experiences are one of the most significant aspects of XR, VR, and AR. The Metaverse takes this to the next level, allowing large groups of individuals to share an immersive first-person experience while maintaining a strong sense of mutual presence. Although the term “metaverse” is often associated with virtual reality, according to Rosenberg, there are two types of metaverses: a “virtual metaverse” in which people are avatars and an “augmented metaverse” in which layers of virtual content are overlaid on the real world and experienced by real people. Figure 2 depicts XR, VR, AR, and the virtual and augmented metaverse, as well as their interaction. Metaverse is the next generation of the Internet that will surround us both graphically and socially. A historical overview of developments regarding metaverse is shown in Figure 3 and different applications of XR in metaverse along with their enabling technologies are illustrated in Figure 4.

### B. Metaverse: An Architectural Overview

The architecture of the metaverse expands from the people’s experiences to the underlying enabling technologies and has

seven layers that include: (1) Experience; (2) Discovery; (3) Creator; (4) Spatial Mapping; (5) Decentralization; (6) Human Interface; and (7) Infrastructure, which is illustrated in Figure 5 and is briefly described below.

- *Layer 1: Experience:* It is the topmost layer in the metaverse, which is mainly concerned with the experiences of the users and it provides different services to them, e.g., games, E-sports, social interactions, events, festivals, shopping, co-working, etc.
- *Layer 2: Discovery:* It is like a push and pull service that introduces people to new experiences in the metaverse such as virtual stores, advertising networks, ratings, social curation avatars, chatbots, etc. It will involve both inbound (i.e., users are actively seeking information regarding experiences) and outbound (i.e., an advertisement that is not explicitly requested by the user). This layer is mainly driven by metaverse service providers and content creators to inform and motivate users regarding new features and services.
- *Layer 3: Creator:* This layer is sometimes also referred to as the creator economy. Like the previous layer, it is mainly driven by the content creator and service providers, who leverage different technologies to create

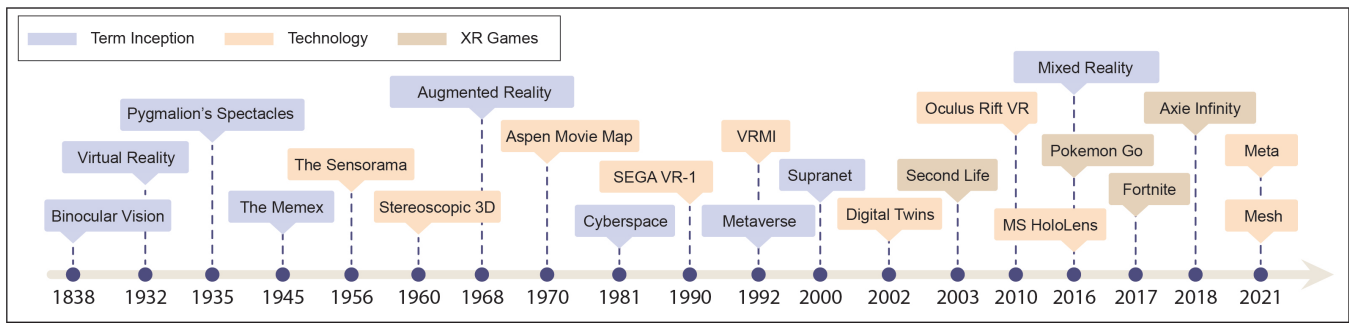


Fig. 3. Historical evolution of technologies needed for metaverse



Fig. 4. Applications of XR in metaverse.

content or experiences for metaverse users such as asset markets, E-commerce, design tools, and workflow.

- **Layer 4: Spatial Mapping:** This layer provides a bridge between the digital world and the physical world and provides immersive experiences to metaverse users. It consists of different technologies like geospatial mapping, object and speech recognition, 3D engines (for enabling animations), VR, AR, XR, multitasking, and integration of user interfaces and heterogeneous sensor data (e.g., from IoT and wearable devices). It can be assumed as the backbone of the creator layer [10].
- **Layer 5: Decentralization:** Decentralization is very crucial in the metaverse and ideally it should not be controlled by a single entity. It provides a scalable ecosystem to developers in terms of providing online capabilities and reliability to the users at the same time. This layer will consist of multiple technologies that include edge computing, blockchain, microservices, and AI agents.
- **Layer 6: Human Interface:** This layer is mainly concerned with the interfacing of the physical world with

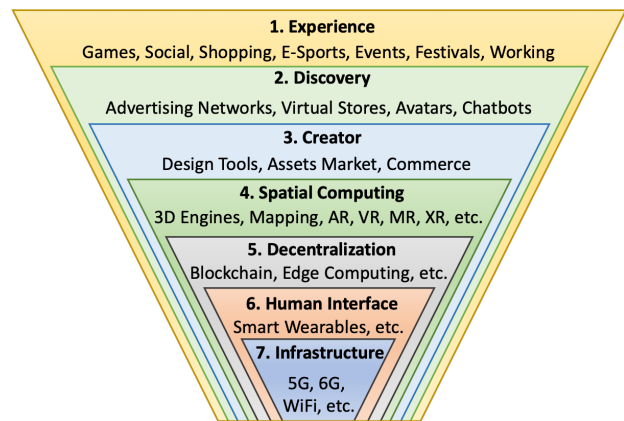


Fig. 5. Illustration of different layers in metaverse.

the digital and from the digital to the physical world. For example, let's consider the example of metaverse services that require data collected from humans using different sensors such as smartwatches, smartphones, smart glasses, wearable IoT devices, biosensors, and head-mounted displays, just to mention a few.

- **Layer 7: Infrastructure:** This layer is responsible for connecting different enabling devices and technologies to the network for content delivery in the metaverse. Different ICT technologies will act as a backbone in the infrastructure layer of the metaverse. For example, 5G/6G-based communication has huge potential to drastically improve the performance of metaverse applications while reducing latency and speeding up content delivery. In addition, this layer will also involve major data processing capabilities like data centers, the cloud, CPUs, GPUs, and even quantum computers.

### C. Metaverse: Applications

Metaverse applications that incorporate different technologies like VR, AR, or XR have various potential applications in education, healthcare, industry, and scientific research, just to name a few. A detailed taxonomy of various potential metaverse applications is illustrated in Figure 6. Metaverse allows moving from text-focused Internet that supports 2D images to a 3D or even a 4D world (in which we may travel in time (forward or backward)). One of the promising



Fig. 6. Various potential AI-XR metaverse applications.

applications will be social VR, which will be the enhanced version of current social media. As metaverse can leverage both VR and AR, numerous applications (e.g., voice recognition, gesture recognition, and speech translation) can benefit [11]. Over the past few years, VR and AR technologies have become very mature and nowadays their equipment is relatively cheap and readily available. This is a long way from the modest beginning of AR and VR, which were ignited in the 1960s by Ivan Sutherland in his pioneering work on the first responsive head-mounted wearable devices, which were admittedly primitive by modern standards.

Modern VR headsets have become accessible (e.g., Facebook’s Oculus Quest 2 is available for \$300) with the price expected to go down as technology continues to advance. There are numerous AR applications such as Heads-Up-Display (HUD) features on modern luxury cars, the use of face filters in apps such as Snapchat, and games such as the addictive Pokémon GO game, where players could “see” Pokémon characters on the street. Modern mobile phones supporting Lidar technology now can support AR with new software development kits emerging such as Google’s AR development platform ARCore, which provides the ability to track motion, understand the environment, and estimate light—three capabilities essential for AR. AR pioneer Louis Rosenberg predicts that within 10 years, most people will be spending more than 2 hours every day in VR, and augmented reality interfaces will replace mobile phones as our primary interface with digital content.

#### D. Role of AI and ML in Metaverse

Different AI techniques including machine learning (ML) and deep learning (DL) have many potential applications in the metaverse (as shown in Figure 7). For example, one of the most fascinating features of the metaverse will be voice-based commands that will utilize different voice recognition and analysis techniques for language processing and understanding human commands. In addition, metaverse will use different ML/DL-based regression and classification for data

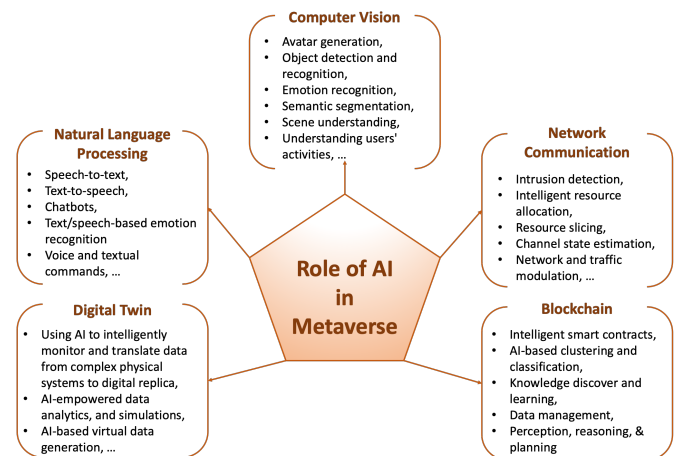


Fig. 7. Applications of AI in metaverse.

management and decision making. Also, to provide immersive experiences to the users it will use different generative models to generate photo-realistic avatars and for 3D reconstruction of objects from 2D images. We discuss the potential applications of ML/DL in the metaverse across five dimensions, which are described next.

1) *Applications of ML in Natural Language Processing:* Natural Language Processing (NLP) consists of different techniques that are used for automatically analyzing and understanding human languages (i.e., text and speech). There are many NLP applications that will be part of AI-XR metaverse applications, e.g., speech-to-text, text-to-speech, chatbots, and text/speech-based emotion recognition are the most prominent features of the metaverse. In particular, NLP techniques will be used for the recognition and understanding of complicated human conversations and commands. A key driving force behind the success of NLP methods is the advancement in ML/DL, with the development of new techniques such as recurrent neural networks (RNN), long-short term memory (LSTM), and transformer networks [12].

2) *Applications of ML in Vision*: Machine vision or computer vision will be a fundamental component of AI-XR metaverse applications. Different computer vision applications will enable various functionalities in metaverse applications, for example, processing visual data from different sensors to infer high-level visual semantics. The major tasks of the visual processing pipeline include understanding users' activities, emotion recognition, object detection, scene understanding, semantic segmentation, avatar generation, etc. In addition to these applications, the metaverse is expected to have AI-empowered quality assessment capabilities, e.g., for satisfying the users' demands about viewing high-resolution videos [6]. In this regard, advanced AI methods can be used to develop quantitative and qualitative benchmarks for visual quality assessment.

3) *Applications of ML in Network Communication*: Metaverse is expected to simultaneously entertain a massive number of users with the metaverse services provisioned mainly through wireless networks. Over the past few years, substantial research attention has been devoted to improving the overall throughput and performance of wireless network communication and the use of different AI techniques is the main driving force behind this innovation [13]. Metaverse will mainly include real-time multimedia services that require a reliable connection, high throughput, and low latency to ensure a seamless user experience. Therefore, it is expected that the metaverse will benefit from 5G and beyond empowered communication. The potential of different AI techniques has already been demonstrated for 5G and 6G, e.g., intelligent resource allocation [14], solving resource slicing problem [15], channel state estimation [16], and network modulation [17], etc.

4) *Applications of ML in Blockchain*: Service providers in the metaverse will provide users with different incentives in terms of digital assets (e.g., coins) for different events, games, and creative activities. The dispersion of such assets requires a transparent way to record and track such transactions. In this regard, smart contracts empowered blockchain technology can be leveraged that allows critical information to be stored on an immutable and impenetrable ledger. The decentralized nature of blockchain imbues it with great potential to address security and privacy issues in metaverse [18]. This potential increases with AI-empowered blockchain applications [19], for example, different AI-based clustering and classification techniques can be used for data analysis stored on blockchain [20]. In addition, different AI techniques can be used for knowledge discovery and learning, efficient data management, perception, reasoning, and planning.

5) *Applications of ML in Digital Twin*: The term digital twin refers to the digital replica (i.e., representation) of real objects. A digital twin is capable of synchronizing regular actions, operations processes, and assets with the real world, e.g., analyzing, monitoring, predicting, and visualizing [21]. The digital twin also acts as a bridge where the actual world and digital world interact with each other through different IoT devices [22]. The digital twin will be one of the most important building sectors of the metaverse that allows users to access and use services in the virtual world while exactly

depicting the real world in a virtual environment. For example, surgeons and medical experts can create a digital replica of a patient to study and understand the involved complexities before performing his surgery.

### III. CHALLENGES IN IMPLEMENTING SECURE AND TRUSTWORTHY AI-XR METAVERSE APPLICATIONS

Despite the significant potential of different AI-XR metaverse applications, there are various challenges related to security, privacy, and lack of trust that can hinder their wide adoption. A few such challenges include privacy breaches, security invasion, user profiling, unfair AI outcomes, etc. These challenges may directly or indirectly put the users' safety at risk and can influence social acceptability [23]. Moreover, as discussed above metaverse is the integration of different modern technologies like AI, blockchain, and 5G/6G, therefore, it is likely that the inherent issues associated with these technologies get translated into the metaverse. In this section, we describe different challenges that can hinder the secure, safe, robust, and trustworthy employment of AI-XR metaverse applications. Specifically, we characterized and discuss these challenges in two dimensions, i.e., challenges associated with the use of AI techniques including ML/DL-based methods, and XR-related challenges in the metaverse. We will start by first discussing AI-related challenges.

#### A. AI Associated Challenges in AI-XR Metaverse

Modern AI techniques that include ML/DL-based models suffer from different vulnerabilities that hinder the smooth, safe, secure, and trustworthy use of these methods in critical applications like healthcare, autonomous vehicles, and AI-XR metaverse applications. Below we briefly discuss various such challenges.

1) *Privacy Issues*: Ensuring the privacy of the end users will be a major challenge in AI-XR metaverse applications. As these applications are designed to monitor and collect users' data at an unprecedented fine-grained level, in a bid to create a replica of the digital world [2], there is a greater danger and risk of privacy breaches [8]. For example, to create an immersive virtual scene in the metaverse, data from different sensors will be collected and analyzed using AI models, e.g., facial expressions, brain wave patterns, hand movements, eye movements, biometric, and speech data [8]. This raises obvious concerns regarding the privacy of users and opens a new horizon for digital crimes [2]. Users' sensitive information including daily routine activities, personal logs, and schedules will be stored on a server, which ultimately becomes a critical privacy challenge in a publicly distributed network. Such data include body movements, voice, reflexes, and even more critical data that include subconscious and unconscious responses such as eye movements and physiological signals. Features such as eye tracking are readily accessible using commercially available products such as the HTC Vive Pro Eye and Pico Neo VR headsets even though the XR expert Louis Rosenberg recommends banning such features and data collection in non-health-related applications for ethical reasons. In general, the data is collected through on-device sensors at the user site,

processed at their devices or nearby local server, and logged as storage in the cloud. Considering the above-mentioned procedure, some malicious attacks can be encountered which are classified as (i) data collection, (ii) data storage, (iii) data usage, and (iv) user profiling. Furthermore, as AR technology depends on the precise localization of users in the physical world, modern smartphones use Lidar sensors and Simultaneous Localization and Mapping (SLAM) algorithms to precisely locate the user in the real world. This opens up the possibility of privacy violations as sensitive information may be exploited for anti-social purposes. Some important concerns related to the privacy of data are described next.

*a) Data Collection:* Generally, the data in AI-XR metaverse applications is collected through users' input either by voice command or textual input, and multi-modal sensors including camera, microphone, textual commands, gesture sensors, and wearable devices. These devices are expected to frequently collect personal information such as daily routine activities, voice and biometric data, and personal preferences including shopping items, TV shows, and preferred food choices. In addition, the private data will be used for the creation of avatars for a digital representation of a real human in the metaverse, which also raises privacy issues. For example, the built-in location sensor in the Oculus headset can be used for tracking users' presence in the real environment with a precise accuracy [8].

*b) Data Storage:* The ultimate aim of developing personalized AI-XR metaverse applications is to aid human beings and ease their daily routine activities. These applications will contain multiple sensing devices which generate a substantial amount of data. Whereas these devices will be resource constrained with small storage units, which leads to the tradeoff between data generation and storage at the edge level. To address these shortcomings, these devices upload their data along with the corresponding logs to online local or global data centers. Though, the data storage tradeoff is resolved by connecting with the servers. However, it also raises privacy concerns regarding the access permissions and data protection of the consumers. Also, if the communication channel is hacked by an attacker then the data can be manipulated to get the intended outcomes. The literature suggests that adversaries can extract information regarding the actual data even if the communication is encrypted and can track the location of the users by realizing different attacks such as advanced inference attacks [24] and differential attacks [25].

*c) Data Usage and Consent:* Continuous data collection through multi-modal sensors including cameras, microphones, and other sensors will be closely involved in the daily routine activities of metaverse users. These sensors collect continuous data regardless of privacy awareness, which is logged over the local or global server(s). Consequently, this procedure raises legal questions regarding the users' consent and the kind of data that is collected and shared. Moreover, metaverse service providers can also utilize this data to optimize their inference models and to make them robust and more personalized. However, it also raises privacy concerns such as the collection, disclosure, and sharing of the data without the explicit or implicit consent of users.

*d) User Profiling:* Similar to the current social media (in which users are considered as the product), everything will be a product alike in the metaverse. Metaverse will act as a meta-platform for different entities (such as users, developers, content creators, businesses, etc.), and this raises questions about data collection and its utilization for user profiling [5]. Also, the provisioning of the metaverse services requires that the users should be uniquely identifiable in the metaverse. For this purpose, VR headsets/glasses or any other such device can likely be used for illegally tracking users in real life [26]. Moreover, such devices can be attacked by malicious actors and can be exploited to track users' real locations for possible digital and real-world crimes. Guzman et al. [3] presented a data-centric perspective to avoid unprecedented privacy challenges related to data collection and its usage in MR.

*2) Lack of Trustworthiness:* According to the definition of Trust, expressed by Lee and See [27], in the perspective of automated systems, "*Trust is an attitude that an agent will help achieve an individual's goals in a situation characterized by uncertainty and vulnerability.*" Metaverse is a data-driven technology in which service providers will use different automated tools to assist human beings as a recommendation engine in various domains, e.g., shopping recommendations, movie recommendations, and even recommendations regarding their health. The efficacy of such recommendation systems is highly dependent on the collection of personalized data for intelligent decision-making, e.g., an AI model is trained using the collected fine-grained health data to suggest more accurate recommendations regarding the well-being of users. Despite the huge potential of such features of AI-XR metaverse applications, it also raises many questions about which parameters have influenced the underlying AI model in producing a particular decision and the trustworthiness of such predictions and decisions. Next, we discuss such challenges across three dimensions (i) Truthful AI, (ii) Transparency, and (iii) Explainable and Interpretable AI.

*a) Need for Truthful AI:* The evolution in technology also brings threats among them. Over the past few years, AI-based personal assistants including Siri, Alexa, Astro, and the like, have managed to gain wide social acceptance and these devices are being used by consumers daily for automating household routine tasks. Currently, misinformation, falsehood, or malfunction in AI-based text or speech analyses and command execution are not considered a matter of concern. However, it is expected that AI-enabled intelligent systems with linguistic capabilities will be a major feature of AI-XR metaverse applications. In this regard, it will be quite challenging to enforce the criterion of truthfulness in AI-based systems to ensure the safe selection of statements and behavior according to the social norms of users while interacting with human society.

*b) Transparency Issues:* In general, practices, developing AI-based systems is about training sophisticated algorithms on large-scale data to learn an efficient and generalizable model that can be deployed in the real-world environment. However, it is well-known fact that the performance of these models is directly proportional to the transparency of training

data, i.e., AI models will perform as well as their training data. Numerous factors such as input data riddled with poorly cleansed, or selection of inherently biased data, underfitting, and overfitting influence the performance of these models and can result in fairness and accountability issues (discussed later in this section). Unlike typical application development, there are no quality assurance tools available to spot bugs and evaluate the bias factor in the training data. For instance, if it was known at which stage, the model is going to infer at a perfect scale, then there would not be a need to perform training on such large-scale data. This process is all about the hit and trial procedure, which is a quite challenging task to identify the right approximations with better data, hyperparameters, and configuration settings.

c) *Explainable and Interpretable AI*: The rapid adoption of based applications in human society has also grown the complexity of the systems, which ultimately requires system understandability to make them legitimate and trustworthy. In a critical human-facing technology such as the AI-XR metaverse, interpretable and explainable AI models are required to answer questions about accountability and transparency of their decisions and outcomes. For example, how the employed model reached the decision, and which factors influenced the models to make that decision [28]. Such questions are particularly important for human-centric applications where the potential impact of AI will be limited if it is not able to provide accurate and transparent AI predictions. Therefore, the key objective of these models is to develop a relationship of trust between human users and AI. However, one of the main challenges in developing explainable methods is the trade-off between achieving the modesty of an algorithm and ensuring the discretion of sensitive user data. In addition, it is also a challenging task to identify the right information while generating a simple yet useful explanation for users. It is worth noting that the terms interpretability and explainability are closely related and are often used interchangeably in machine learning literature, however, these terms are different in practice. Interpretability of the models is defined as the extent to which its outcomes are predictable, i.e., for a given change in the input or model parameter(s), the interpretability enables us to predict the respective change in its output. On the counter side, explainability deals with the explanation of internal processes of the ML/DL models in a human-understandable way.

3) *ML Security Issues*: Despite the state-of-the-art performance of modern AI techniques including DL-based systems, it has been shown that these models are highly vulnerable to carefully crafted adversarial perturbations known as adversarial ML attacks [29]. The threat of these attacks has been already demonstrated for many critical applications like healthcare, autonomous vehicles, etc. On a similar note, AI-XR metaverse applications are essentially critical as they involve humans, and ensuring their safety from any harm is profoundly important. On the counter side, the existence of these challenges raises many concerns about the safety, security, and robustness of AI-based metaverse applications thus hindering their practical deployment. As it is equally important that any AI-based should be equally trusted by all

stakeholders involve service providers, developers, and end users. These challenges are detailed later (Section IV).

4) *Lack of Fairness and Accountability*: Modern AI methods like advanced DL models lack fairness and accountability in their decisions [30]. On the other hand, such questions are particularly important for critical applications like AI-XR, in which the model's decisions can have life-threatening consequences for the end users. Moreover, AI models are developed using training data, which will be mainly collected from human users in AI-XR metaverse applications for providing immersive experiences. Humans possess certain biases that will be readily reflected in the data they generate, and when this data is used for training AI models, the data bias will be directly translated into the developed AI-based system. As a result, the model will be biased towards certain samples that contain certain features (bias), and its decision will not be fair. On a similar note, the critical nature of AI-XR metaverse applications demands accountable decisions. Consequently, data bias if remained unaddressed can ultimately lead to unintended consequences [31].

5) *Identity Theft and Authentication Attacks*: Users'/avatars' identities in the metaverse can be stolen or impersonated illegally leading to authentication and access control issues in the interconnected virtual worlds. Identity theft in the metaverse will be more dangerous than traditional attacks. The identity of a user once stolen will reveal everything about that person's digital assets, avatars, and social relationships. The attackers can exploit different vulnerable VR gadgets and other service authentication loopholes to realize identity theft attacks and can steal the victim's secret keys of digital assets and bank details. It has been reported that about 17 users in the OpenSea NFT marketplace were hacked through a phishing attack and flaws in the smart contract that resulted in a loss of \$1.7 million.<sup>1</sup>

Metaverse will leverage different biometrics and password-driven technologies to authenticate users and their avatars in the virtual worlds. The attacker can evade such authentication systems to impersonate real users' identities to get control of the whole virtual world. Evading AI-based biometric systems has become easier with the advancements in adversarial ML research. Therefore, AI-empowered speech and face recognition-based biometric systems can be easily attacked to realize impersonation attacks. Once the attacker has the access to the metaverse it can exploit the data generated by the victim's devices to deceive him, committing a crime in the virtual space. On other hand, the exposure of biological data when used for authentication purposes can also lead to severe consequences [32]. Moreover, the authentication of social friends of a user using their avatars is much more challenging in the metaverse as compared to real-world identity authentication. In this regard, facial data, voice, and videos can be used to develop an AI-based avatar authentication system, however, the unsolved inherent issues of AI can still hinder its practicality.

<sup>1</sup><https://threatpost.com/nft-investors-lose-1-7m-in-opensea-phishing-attack/178558/>



6) *The Bias Problem*: Bias refers to a model making certain unethical assumptions about the data. Human bias along with its many aspects has been studied by researchers in many disciplines including law, psychology, and so forth. In [33], bias is defined as ‘the prejudice or inclination of a decision made by an AI system which is in a way considered to be unfair for or against one person or group’. Bias in recommendation systems, advertising algorithms, facial recognition systems, and risk assessment tools has been widely studied in recent years. In metaverse applications, data will be collected from a heterogeneous group of people and sources having their own characteristics, stereotypes, and behaviors, which introduces different biases in the collected data. In [34]–[36], the authors discuss different kinds of bias based on the sources and the types of bias. On the base of sources, these biases have been divided into further categories: biases caused by data, biases caused by algorithms, and biases caused by user interaction.

### B. XR-related Challenges in Metaverse

AI-XR metaverse applications are essentially human-centric and ensuring the security, privacy, security, and robustness of such applications is of utmost importance. It has been envisioned that an entirely new form of digital media will emerge from the use of VR and AR in the modern metaverse (TV, print media, and the web). In recent years, there has been immense discussion regarding the concerns about surveillance capitalism, which is happening on the Internet in different applications. Many large tech organizations providing Internet services like Facebook, Google, Microsoft, and Amazon collect large of amount data related to the surveillance of their users, which is then used to satisfy the needs of advertisers [37]. The pioneers of VR and AR such as Jaron Lanier [38] and Louis Rosenberg<sup>2</sup> have predicted that the concerns about surveillance are expected to rise in Metaverse. For example, it has been shown how reconfiguring AR in Pokémon Go (an AR mobile game) drew unexpected audiences to museums and public spaces like trains to fill in the space thus creating a form of virtual trespassing. It has been reported that people were putting their lives in danger to pursue virtual characters. This highlights that safety concerns may arise when such immersive technologies are engineered for gaming and experiences. Below we discuss the key challenges that are hindering metaverse applications in general and we will later discuss the specific challenges that arise with the use of different AI techniques in metaverse applications (Section IV).

1) *Safety Issues*: There are various concerns regarding the mental and physical safety of metaverse users. There are several reported incidents of digital harassment, theft, and bullying in XR applications [39]. The report on “Immersive and Addictive Technologies” highlights rampant incidents of sexual harassment, cyberbullying, and grooming online [40]. Ensuring the safety of users is a major challenge for AI-XR metaverse applications because of the fact that such incidents have real damage and harm to users despite being experienced in the virtual world. The avatars generated using

recent advancements in AI techniques, in particular, generative models can appear more realistic in AI-XR metaverse applications and can engage users in promotional conversation thus providing a false sense of a real human behind the avatar. The avatars in such a promotional are fueled with more personalized data (such as your vitals, emotions, expressions, etc.) to look more realistic. Also, these sales avatars can pitch products to you more persuasively than any real salesman or even a recommendation system due to their access to rich cyber-physical data about you. The research in deep fake technology and photorealistic avatars is already on the stage where computer-generated content is indistinguishable from the real. Such advancements can be leveraged to realize an adversarial attack on AI-XR metaverse applications to get the intended behavior and outcomes.

2) *Potential Antisocial Aspects*: There are various opinions regarding the antisocial aspects of AI-XR metaverse applications, many people think that introducing AI-XR metaverse may detract the users (humans) from their purposes and may have a somatic effect. In the literature, it has been shown that extended times online can result in users demonstrating post-VR sadness and detachment from reality. For instance, Aldous Huxley in 1932 wrote in his social science dystopian fiction novel that using technology can lead to self-inflicted harm that can lead people to be diverted from their higher priorities and become more prone to being influenced by other interests. As a result of such a quest for technological utopia, the human psyche and society as a whole are greatly afflicted. Social critics have long argued that various digital media, such as television, the Internet, and the Web, make people docile and less connected to the real world. For example, Jerry Mander in 1978 wrote in his book, “Four Arguments for the Elimination of Television” that TV removes the sense of reality from people, promotes capitalism, TV can be used as a scapegoat, and all these three factors work together negatively. The modern technological disruptions including the web and social networking services have created a filter bubble detaching people from the real world and the truth. Due to these reasons, the current era is also referred to as a post-truth era [41]. We may reasonably expect that alienation from the real world will exacerbate with the increasing adoption of VR, AR, and AI-XR metaverse applications, which aim at changing the human perspective of the world. This argument can be supported by the fact that in 2018, the World Health Organization formally included “gaming disorder” in its International Classification of Diseases following research that shows that technology can promote addictive behavior in people. Moreover, the literature focused on analyzing the social implications of metaverse argues to understand and identify potential psychological problems that can arise in metaverse [42].

3) *Ethical Aspects*: Any technological intervention involving humans suffers from some serious ethical issues, especially the one that contains intelligence. The Institute of Electrical and Electronics Engineers (IEEE) has recently published a report on Ethically Aligned Design that mainly focuses on the Ethics of Autonomous and Intelligent Systems [43]. This report emphasizes the need of developing ethical autonomous and intelligent systems (A/IS) that promotes human wellbeing

<sup>2</sup><https://bigthink.com/the-future/metaverse-augmented-reality-danger/>

and protects human rights through transparent and accountable A/IS and the prevention of the misuse of AI. This report is the collective effort of hundreds of researchers having diverse backgrounds and expertise in important areas like governance, technology, civil society, and policy-making. This report has a dedicated section on XR and interestingly, IEEE also has a Global Initiative on Ethics of XR.<sup>3</sup> In this report, various ethical issues related to XR have been highlighted including users' preference for virtual life over the real world and complete disengagement with society. In addition, the reports conclude with the following remark regarding XR: *"The nature of XR environments fosters unique legal and ethical challenges that can directly affect users' privacy, identity, and rights. Society will need to rethink notions of privacy, accessibility, and governance across public and private spaces. New laws or regulations regarding data ownership, free use, universal access, and adaptive accessibility within XR environments may need to be developed."*

4) *Regulatory Challenges:* The big tech companies have resisted regulation decriing the fact that regulation will slow down innovation. However, there are various ethics researchers and social scientists who are arguing for much greater regulation to ensure that consumer rights are protected. In this regard, the EU General Data Protection Regulation (GDPR) has paved the way for many countries and regions attempting to develop similar regulatory laws to protect Internet users. These regulations mainly emphasize the importance of the non-profiling of users (i.e., limiting the storage of tracking data) and for better transparency (i.e., online services and applications should specify why and what information is being stored). On a similar note, AI-XR metaverse applications are subject to the requirement of being transparent in terms of data collection and utilization and should also be subject to the informed consent of users. Also, the development of such applications requires thoughtful deliberation from a regulatory perspective. For instance, is worth considering banning non-medical applications to collect vital biomedical statistics due to the high risks of being exploited maliciously. To mitigate the risk of users being manipulated deceptively, metaverse operators may be bounded to transparently declare the staging of virtual products and experiences in the metaverse. Rosenberg, one of the pioneers of VR and AR, has already started to argue about the need for regulation for metaverse applications [44]. For instance, he suggested leveraging the arguments regarding the regulation of social media for developing a legal and philosophical basis for metaverse regulation. As the metaverse can be deemed as an evolutionary expansion of similar services. Rosenberg argued that the only solution to eliminate ethical and privacy-related concerns associated with metaverse is to shift from an advertising-based to a subscription-based business model in which users pay a subscription fee for accessing the metaverse platform. This eliminates the service providers' need to monitor their user base to a greater extent, however, this is not a feasible solution as it is difficult to say whether or not users will pay for a safer metaverse.

#### IV. ANALYZING SECURITY AND TRUSTWORTHINESS ASPECTS OF AI-XR

In this section, we will discuss the challenges associated with the use of different AI techniques (in particular, ML/DL-based models) that hinder the safe, secure, and trustworthy deployment of these methods in metaverse applications. We start by first providing a broad overview of AI security in the metaverse.

##### A. Security of ML in Metaverse

The impact of metaverse applications will be social and economic and these applications will be more susceptible to undesirable adversarial action(s). AI will be the fundamental driving force behind the success of the metaverse, there are numerous applications of AI in different layers of the metaverse. On the other hand, the use of AI algorithms in AI-XR metaverse applications also opens them up to different adversarial attacks. In Figure 8, we highlight the threat of different security and privacy attacks that can be realized in different applications in almost every layer of the metaverse. The figure also highlights that there are various common ML security issues and attack surfaces that get shared across the architectural landscape of the metaverse across different AI applications at each level. In this section, we discuss different AI-associated security and privacy attacks on AI-XR metaverse applications.

##### B. Potential Attacks on AI-based Metaverse Applications

The threat of adversarial ML attacks has already been shown to be successful in compromising the integrity of AI techniques in many critical tasks, e.g., connected and autonomous vehicles [45], computer vision [46], and healthcare [30], just to name a few. Furthermore, AI-XR algorithms could be biased either due to data imbalance or adversarial subversion. Many of the ethical dilemmas and social harms such as distraction, narcissism, disinformation, outrage, and polarization stem from the economic model of surveillance capitalism in which service providers give the customers everything and anything that makes the company money. In this way, these companies pander to the base animal desires of people and exploit their cognitive biases effectively downgrading humans and manipulating them for ulterior selfish purposes.

In the adversarial ML literature, an adversarial example is defined as the input to the deployed AI model crafted by an adversary by introducing imperceptible noise into the legitimate sample to get the intended outcomes. In general, there are two types of adversarial ML attacks: (1) poisoning attacks that aim at altering the training process of the AI model; and (2) evasion attacks that are focused on evading the deployed AI model by making inferences (they are also known as inference time attacks). In poisoning attacks, the adversary mainly modifies the training data to tamper with the learning of the AI model [47]. In contrast, test data is manipulated in evasion attacks to get the desired predictions from the model [48]. Recent works have shown that AI models are vulnerable to attacks at both training and inference

<sup>3</sup><https://standards.ieee.org/industry-connections/ethics-extended-reality/>

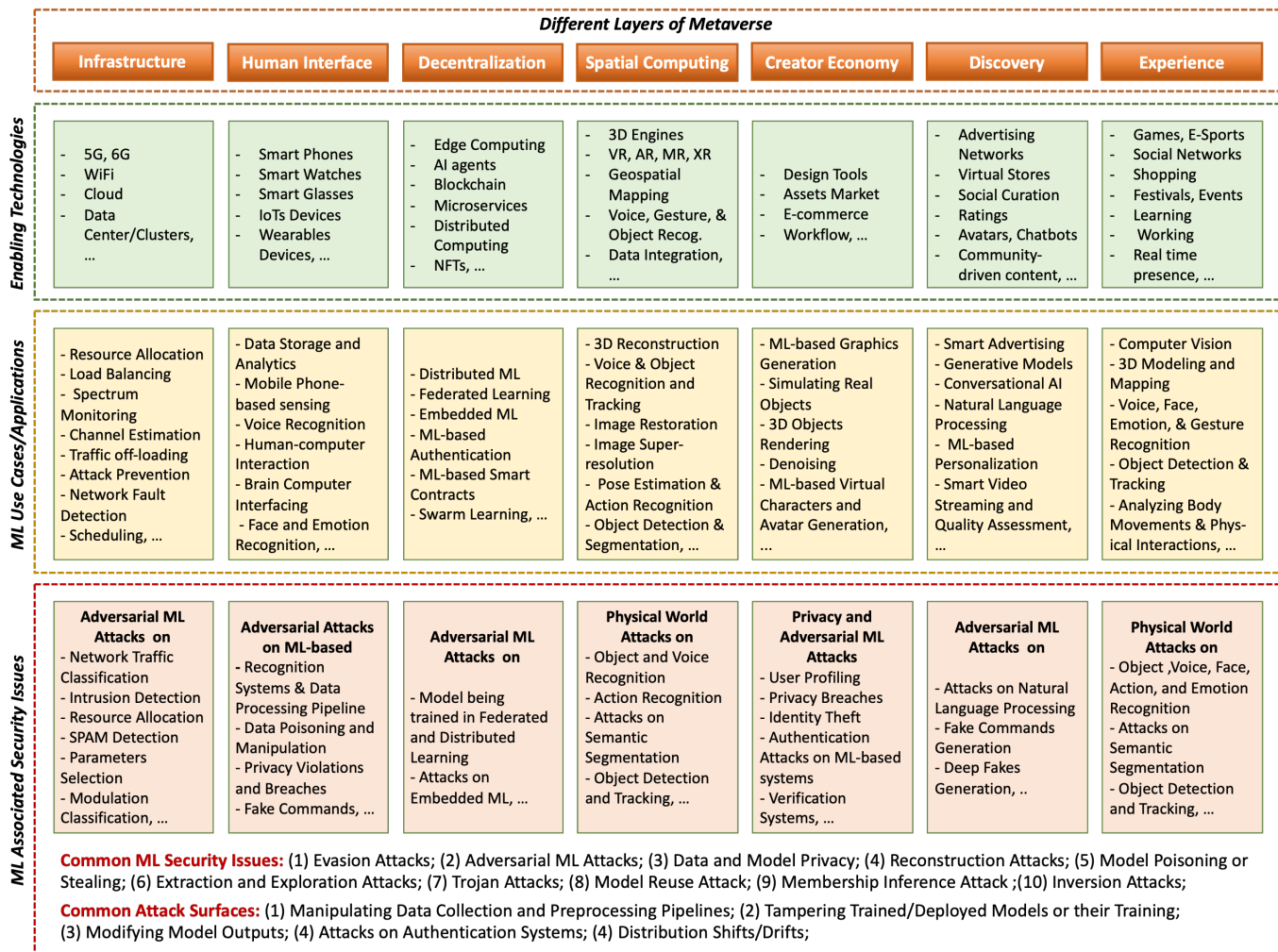


Fig. 8. Overview of ML security in Metaverse.

stages [49]. Training stage attacks typically corrupt a small subset (typically  $\sim 1\%$ ) of the training data samples to achieve malicious goals during AI model training [50], [51]. On the other hand, the inference stage attacks cause a trained model to misbehave on adversarially crafted test inputs [52], [53].

Attacks on AI models are generally carried out by first defining a threat model. A threat model is a set of assumptions regarding the attackers' abilities to access and affect a typical AI model training pipeline. Broadly, there are two main threat models—the poisoning threat model (i.e., realizing poisoning attacks), and the adversarial threat model (realizing evasion attacks). A poisoning threat model assumes an attacker who can control a small set of the training dataset to adversely affect the training of the model. An adversarial threat model assumes an attacker who can access and, to a certain extent, perturb the inputs to an already trained AI model. In the following, we highlight major security threats associated with the use of AI techniques in metaverse applications that include computer vision, natural language processing (MLP), network communication, authentication, and recognition systems.

1) *ML Associated Security Issues in Computer Vision:* Computer vision is one of the central building blocks in the

foundation of the metaverse. In recent years, DL algorithms have enabled major advances in computer vision ranging from image classification [54], [55] to scene understanding [56], [57] and generating realistic images [58]. However, the discovery of the adversarial vulnerabilities of DL-based image processing models by Szegedy et al. [29] sparked a growing concern regarding the reliability and security of these deep models [59]–[62]. Numerous works have analyzed these adversarial vulnerabilities in greater depth under different threat models [49]. In general, adversarial attacks work by optimizing the perturbation,  $\Delta x$ , to an input image,  $x$ , such that the output of the model,  $\mathcal{F}$ , is significantly changed, maximize  $\|\mathcal{F}(x) - \mathcal{F}(x + \Delta x)\|$ .  $\Delta x$  is typically optimized based on the gradients which are either computed directly (white-box scenarios) or estimated by introducing random noise (black-box scenarios). Summary of various adversarial ML attacks on different computer vision applications can be seen in Table II.

2) *ML Associated Security Issues in NLP:* Similar to the vulnerabilities of ML models for vision applications, the literature demonstrates that the ML methods for modeling NLP tasks are also vulnerable to malicious attacks, at both the

TABLE II  
SUMMARY OF DIFFERENT ADVERSARIAL ATTACKS ON VARIOUS COMPUTER VISION APPLICATIONS (THAT ARE EXPECTED TO BE POTENTIAL METAVERSE APPLICATIONS).

Application	Authors	Methodology	Datasets	Before → After
Face Authentication	Goswami et al. [63]	Studied how different architectures affect adversarial vulnerabilities.	MEDS, PaSC	89.3% → 41.6%
	Sharif et al. [64]	Developed adversarial glasses to fool face recognition systems.	Celebrity Face	98.95% → 0%
	Shen et al. [65]	Developed black-box attack for face recognition systems using visible light.	CusFace, LFW	100% → 7.9%
	Chatzikiriakidis et al. [66]	Perturbed facial images to fool automatic face recognition to secure a person’s identity.	CelebA	97.8% → 4%
	Dabouei et al. [67]	Studied the vulnerability of face recognition systems against geometrically perturbed faces.	VGGFace2	100% → 0.14%
	Zhong et al. [68]	Used dropout and feature-level attacks to improve the transferability of adversarial inputs.	VGGFace2	100% → 3.24%
	Dong et al. [69]	Used evolutionary algorithm to find adversarial inputs against the models’ decisions.	LFW	100% → 0%
	Wenger et al. [70]	Proposed improved physically-realizable attack against face recognition.	VGGFace	100% → 10%
	Ali et al. [51]	Proposed multi-trigger backdoor attack against backdoor defenses.	Celebrity Face	88% → 8%
	Xue et al. [71]	Exploit hidden facial features as triggers of the backdoor attack.	VGGFace	100% → 0.02%
Object Detection	Zhang et al. [72]	proposed generalizable contextual adversarial perturbations against object detectors.	PascalVOC, COCO	78.8% → 1.6%
	Lee et al. [73]	Showed that non-overlapping physical patches can fool object detectors.	COCO	55.4% → 0.05%
	Xie et al. [74]	Proposed multi-targeted adversarial attacks to fool object detectors.	PascalVOC	72.07% → 3.36%
	Xie et al. [74]	Showed that multi-targeted adversarial attacks against object detectors are transferable.	PascalVOC	54.87% → 37.9%
	Wei et al. [75]	Utilized generative methods to efficiently obtain transferable adversarial inputs.	PascalVOC	43% → 3%
	Wang et al. [76]	Utilized position and label information to attack black-box object detectors.	PascalVOC	100% → 16%
	Wu et al. [77]	Leveraged natural rotations to insert a backdoor into the object detectors.	PascalVOC	89.5% → 4.45%
3D-Object Modelling	Wang et al. [78]	Optimally generates adversarial perturbations against 3D-Object detectors.	KITTI	84% → 0%
	Xiang et al. [79]	Generated 3D adversarial point clouds against PointNet model.	ModelNet40	93% → 0%
	Hamdi et al. [80]	Exploited an auto-encoder to generate transferable 3D adversarial perturbations to point cloud.	ModelNet40	93% → 5%
	Meloni et al. [81]	Used off-the-shelf 3D surrogates to transfer attack on 3D object models.	N/A	100% → 0%
	Li et al. [82]	Proposed a novel formulation to develop backdoor triggers against 3D point cloud models.	ShapeNetPart	98.4% → 0.5%
Semantic Segmentation	Arnab et al. [83]	Performed an in-depth study of adversarial vulnerabilities of semantic segmentation models.	Cityscapes	77.1% → 19.3%
	Xie et al. [74]	Proposed multi-targeted adversarial attacks to fool semantic segmentation models.	PascalVOC	72.07% → 3.36%
	Hendrik et al. [84]	Analyzed universal adversarial perturbation to fool a segmentation model for any input.	Citscapes	64.8% → 12.9%
	Li et al. [85]	Poisoned the segmentation models using object-level target class and semantic triggers.	ADE20K	37.7% → 25.2%
	Feng et al. [86]	Proposed frequency-injection backdoor attack against medical image segmentation tasks.	KiTS-19	54.5% → 21.1%

training and the inference stages of a typical ML pipeline [51], [52]. Below we discuss such attacks.

*Poisoning and Trojaning Attacks:* Poisoning attacks and trojaning (also known as backdoor) attacks are the most widely known training stage attacks in NLP. Poisoning attacks aim to tamper with the training of the model so that it is unable to perform satisfactorily on the test inputs [47], [87], [88]. Trojaning attacks aim to insert a trojan—typically characterized by a specific pattern of words known only to the attacker in the input sequence—into a model such that the model behaves normally on natural test inputs, but malfunctions as desired by the attacker (through the use of of the trojan pattern of

words [50], [89]). Similar to the case with the computer vision applications, the trojaning attacks, being more difficult to be detected as compared to the poisoning attacks, pose a greater threat to NLP applications in the metaverse.

*Adversarial Attacks.* Although adversarial examples have been extensively studied in computer vision, they have received significantly limited attention in NLP tasks mainly due to the discrete input search space—minimal adversarial perturbations in the input are no longer feasible in NLP [28], [52]. Recently, however, there have been numerous works highlighting the adversarial vulnerabilities of the NLP-based ML models. Notable adversarial attacks include Text-bugger, Text-fooler, PWWS,

TABLE III  
SUMMARY OF DIFFERENT ADVERSARIAL ML ATTACKS ON VARIOUS NLP APPLICATIONS.

Application	Authors	Methodology	Datasets	Before → After
Language to Language Modelling	Zhand et al. [90]	propose word saliency speedup local search method to attack translation machines.	NIST (MT)	92.48% Degradation
	Boucher et al. [91]	Uses invisible characters, homoglyphs and deletion control characters to fool the model	WMT14	37% → 1%
Fake-news Detection	Li et al. [92]	exploit BERT-MLM to fool a fine-tuned BERT model by generating coherent perturbations.	AG news	94.2% → 10.6%
	Ali et al. [53]	propose an adaptive adversarial attack to generate perturbations against statistical defenses.	Kaggle Fake news	95% → 0%
	Jin et al. [93]	identify key contributing words and replace them with synonyms while retaining the coherence.	Kaggle Fake news	96.7% → 15.9%
	Zellers et al. [94]	Present a generative model, Grover, to generate fake-news that fools fake-news detectors.	Not applicable	95% → 67%
	Pan et al. [95]	Exploits linguistic styles as triggers to backdoor an NLP model.	COVID	95.1% → 6.7%
Toxicity/Sentiment Classification	Garg et al. [96]	exploit BERT-MLM to generate adversarial perturbations that are coherent with the context.	Amazon	96% → 11%
	Boucher et al. [91]	Uses invisible characters, homoglyphs and deletion control characters to fool the model.	Wikipedia Detox	95% → 19.5%
	Ebrahimi et al. [97]	Leverage atomic flip operation to swap tokens to fool NLP fake news classifiers.	AG news	92.35% → 27.7%
	Li et al. [98]	Exploit model gradients to find and perturb the most positively contributing words.	IMDB	90.7% → 0%
	Li et al. [92]	exploit BERT-MLM to fool a fine-tuned BERT model by generating coherent perturbations.	IMDB	90.9% → 11.4%
	Jin et al. [93]	identify key contributing words and replace them with synonyms while retaining the coherence.	Yelp	93.8% → 1.1%
	Chen et al. [99]	Use char- and word-level triggers to backdoor NLP sentiment classifiers.	SST-5	55% → 0%
	Irtiza et al. [100]	propose a context-aware hidden trigger backdoor attack against NLP classifiers.	IMDB	84.5% → 2.48%

and BERT Adversarial Example (BAE).

Adversarial attacks against NLP models generally follow three major steps—evaluation, perturbation, and selection—to achieve some adversarial goal—for example, targeted or untargeted misclassification—under a predefined threat model. Consider, for example, an input sequence  $X = \{x_1, x_2, \dots, x_i, \dots, x_n\}$  correctly classified by an NLP model,  $\mathcal{F}$ , in class,  $\mathcal{F}(X) = y \in \mathbb{R}^M$ . At the evaluation stage, the attacker uses some *impact scoring function*, to compute a set,  $I_x$ , representing the impact of each word over the output. At the perturbation stage, the attacker repeatedly perturbs the most impactful words in  $I_x$  using some pre-defined *perturbation mechanism* such that the semantic and contextual value of  $X$  remains preserved. At the selection stage, the most optimal perturbation is selected. Table III provides a summary of various adversarial ML attacks on different NLP applications.

3) *ML Associated Security Issues in Networking*: AI-XR metaverse applications will provide ubiquitous connectivity to a massive number of users over wireless networks. Over the last few years, many AI-based algorithms have been developed to improve the performance of wireless communication and networking systems that will be used in different layers of network architecture [13]. The use of AI in wireless communication empowers wireless devices to perform many important intelligent functions such as network composition, analyzing traffic patterns, managing content requests, analyzing wireless channel dynamics, etc. Moreover, AI-based algorithms have been used for optimizing different network constraints like high throughput and low latency for different multimedia applications. A prominent use case is to leverage intelligent proactive load management in 5G and 6G communication networks and predictive data analytics to improve network operations. Despite the significant potential of using various AI algorithms for different optimizing applications in wireless networks, recent studies have highlighted that AI application

is highly susceptible to adversarial ML attacks. For instance, Usama et al. [101] used a generative adversarial network (GAN) for realizing adversarial attacks on network intrusion detection. The threat of adversarial ML attacks on network traffic classification is demonstrated in [102] and for cognitive self-driving networks is presented in [61], [62]. Similarly, the threat of adversarial ML for 5G networks is analyzed in [103]. Summary of various adversarial ML attacks on network applications is presented in Table IV. We refer interested readers to a detailed survey highlighting the threat of adversarial ML in network security [104].

In addition to the above motioned adversarial vulnerabilities associated with the use of AI techniques in many network applications, some other critical network-related issues can hinder the smooth operation of metaverse at a global level. For instance, centralized network architecture provides flexibility in terms of cost saving, simplicity, and ease in performing different operations. On the other hand, such architectures are more prone to a single point of failure (SPoF) and distributed denial of service (DDoS) attacks [108]. For example, if a powerful attacker gets control of the network, it may lead to severe challenges like SPoF and DDoS. To address such issues, the literature suggests leveraging decentralized network architecture [109]. In addition, decentralization will potentially amplify the transparency and trust of users in exchanging their virtual belongings (like digital assets and virtual currencies) among each other and across different virtual worlds in the metaverse. However, many issues arise with the use of decentralized approaches, e.g., reaching a consensus on an ambiguous operation among the huge number of entities in a dynamic metaverse. *Distributed Denial of Service (DDoS)*: Metaverse will include a massive number of IoT devices, which can be compromised by an attacker to form a botnet to realize DDoS attacks [110]. *Sybil Attacks*: In a Sybil attack, the adversary pretends to have fake (or manipulated) identities

TABLE IV  
SUMMARY OF DIFFERENT ADVERSARIAL ML ATTACKS ON DIFFERENT NETWORKING APPLICATIONS.

Application	Authors	Methodology	Datasets	Before → After
Intrusion Detection	Usama et al. [101]	Exploited GAN to craft adversarial examples to evade intrusion detection model.	KDD99	89.12% → 56.55%
	Aiken et al. [105]	Perturbed a few features to evade four ML classifiers trained for detecting DDoS attacks.	KDD99	100% → 0%
Network Traffic Classification	Usama et al. [102]	Crafted adversarial examples using mutual information in black-box settings.	UNB-CIC Tor Data	96% → 77%
Modulation Classification	Usama et al. [103]	Used C&W attack to evade traffic modulation classifier.	RML2016.10a	85% → 15%
	Sadeghi et al. [106]	Realized white-box and black-box attacks on VT-CNN model using a PCA-based perturbations.	GNU Radio	75% → 38%
Network Modulation	Usama et al. [103]	Realized black-box attack on channel autoencoder on unsupervised and DRL models.	RML2016.10a	95% → 80%
Malware Classification	Usama et al. [61]	Realized three SOTA adversarial ML attacks, i.e., FGSM, BIM, and JSMA.	Malware Image Data [62]	98.39% → 1.87%
Abnormal KPI Detection	Usama et al. [62]	Leveraged two SOTA attacks to evade ML-based abnormal KPI detection classifiers.	LTE network data	98.8% → 13.7%
Channel State Estimation	Sagduyu et al. [107]	Realized three attacks: spectrum poisoning, jamming, and priority violation.	Not Articulated	95.58% → 23.12%

of legitimate users or devices. Using such stolen identities he can take over the network.

4) *Security Issues in Cloud-hosted ML Models*: Outsourcing the training of ML/DL models to third-party services that offer powerful computational resources on the cloud is prevalent nowadays. These services allow ML developers to upload their data and models for training over their cloud platforms. It is expected that such services will be featured in AI-XR metaverse applications, as they provide the flexibility of developing AI models using sufficiently large training datasets while reducing the cost and time. However, the literature demonstrates that such services are vulnerable to variety of attacks such as backdoor attacks [111], exploration attacks [112], model inversion [113] and model extraction attacks [114], etc. More details about various attacks and defenses for cloud-hosted ML models can be found in [115]. Visual illustration of adversarial ML attacks on different potential applications in the AI-XR metaverse is presented in Figure 9.

### C. Attacks on VR

The literature highlights that VR systems are vulnerable to adversarial attacks. For instance, Casey et al. [120] demonstrated that humans in VR systems can be controlled like joysticks—thus providing the adversary the ability to control the movement of VR user without his consent or getting into his knowledge. Moreover, the literature suggests that both security and privacy attacks can be realized on VR/XR systems [121]. Therefore, developing secure and robust AI-XR metaverse systems is crucial to the widespread adoption of metaverse applications that are not vulnerable to adversarial attacks or are capable to withstand such attacks and mitigating their impact.

### D. Analyzing Implications of ML Security, Privacy, and Trust Issues: An AI-XR Case Study

In this section, we present an ML/DL-based pipeline for a potential AI-XR metaverse application use case. We then analyzed various challenges and threats that can arise at each development stage. The pipeline is developed while considering a general metaverse application—a virtual conference, in which the participants are remotely connected from different

places (the pipeline is presented in Figure 10). A unique avatar is representing each participant while each avatar is expected to reflect real-time voice, facial expressions, and gestures. The voice of each participant is translated into the native language of all the participants along with generating the transcription. The pipeline depicts different ML/DL-empowered tasks: (1) *3D/4D Visual Reconstruction*—responsible for generating photo-realistic avatars; (2) *3D Visual Mapping*—to reflect real-time multi-modal expressions (i.e., audio, facial, and gestures, etc.); (3) *Speech Recognition and Synthesis*—to interpret and translate the voice of recipient into other languages; and (4) *Speech-to-Text Synthesis*—to generate the transcription of audio conversations of all members.

As depicted in Figure 10, data acquisition is performed by collecting raw audio input through a microphone for NLP, whereas, depth cameras and laser scanners are used for 3D/4D visual reconstruction and mapping. In the next step, acquired data is pre-processed through several techniques including data denoising, deblurring, silence removal, etc. The processed data is then labeled for the training of ML/DL models in a supervised/semi-supervised learning fashion. After successful data preparation, 3D visual mesh construction and segmentation models are trained to perform 3D avatar reconstruction. On the other hand, acoustic and language models along with neural vocoders are trained to perform the multilingual translation and transcription tasks. Although in the literature, these pipelines have demonstrated significant performance in 3D reconstruction and NLP tasks, however, this pipeline is highly exposed to the various privacy and security attacks at each stage of the development pipeline (as shown in Figure 10).

### V. TOWARDS DEVELOPING SECURE AND TRUSTWORTHY AI-XR

The development of secure, safe, and trustworthy AI-XR metaverse applications is fundamentally very important, in this section we will discuss different potential solutions that can be leveraged to address challenges associated with the use of AI in particular and for the overall system in general. An abstraction of different techniques that can be leveraged to

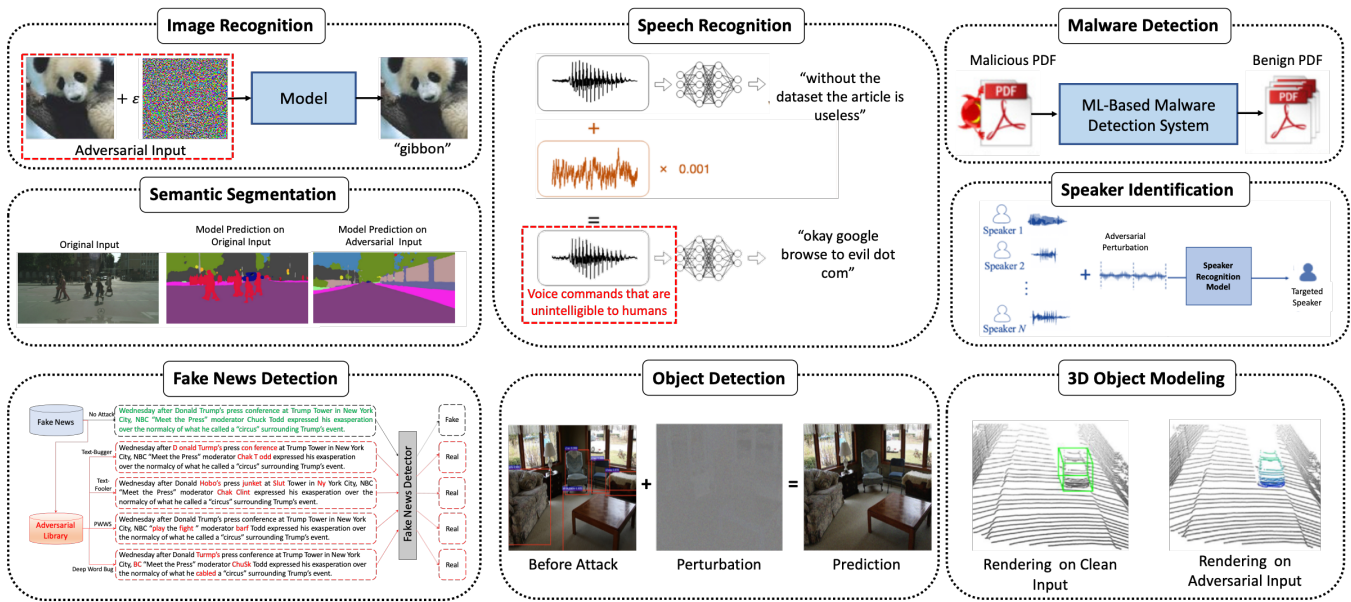


Fig. 9. Illustration of adversarial ML attacks on different potential applications in AI-XR metaverse. Individual figure references: Image Recognition [29]; Speech Recognition [116]; Malware Detection [117]; Semantic Segmentation [118]; Speaker Identification [119]; Fake News Detection [52]; Object Detection [76]; and 3D Object Modeling [78].

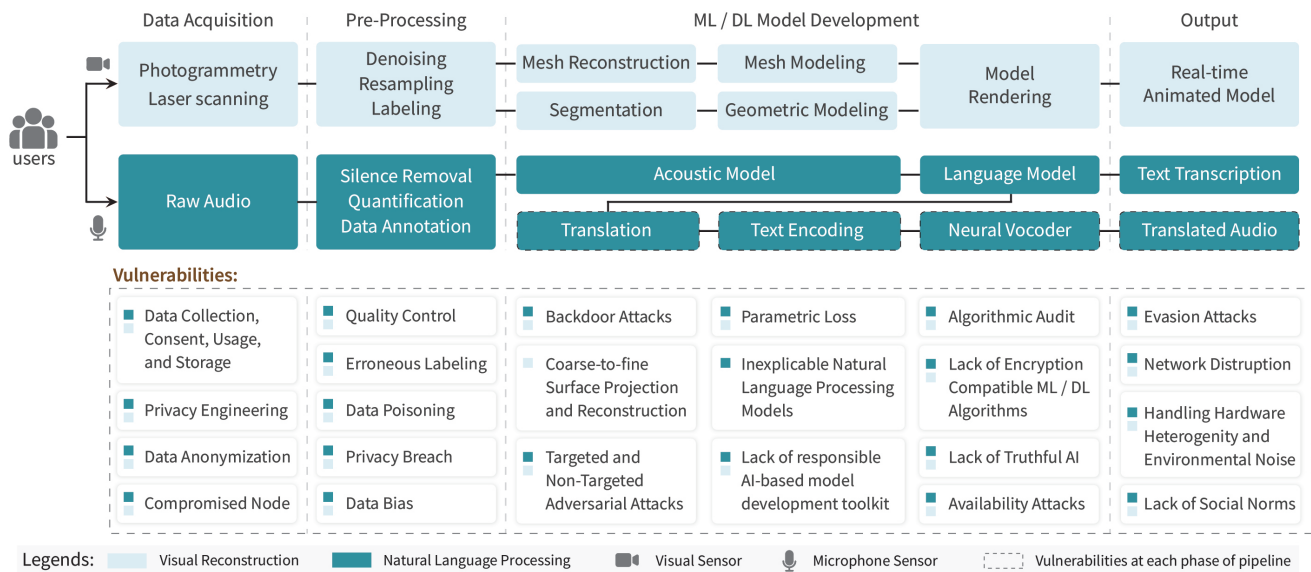


Fig. 10. A prospective pipeline for developing AI-XR metaverse applications for multi-lingual communications, involving various security challenges at each stage.

address the ML-associated issues is shown in Figure 11 and these methods are described next.

### A. Solutions for Privacy Protection in AI-XR

In the literature, privacy-preserving techniques are broadly categorized into three classes: (i) cryptographic techniques, (ii) differential privacy, and (iii) federated and distributed ML. These techniques are briefly discussed below.

1) *Cryptographic Techniques*: Cryptography refers to a practice of methodologies, aiming to construct and analyze communication protocols to ensure secure communication while achieving data integrity, authentication, non-repudiation,

and data confidentiality. Generally, there are two common types of encryption methods: (i) symmetric encryption, and (ii) asymmetric encryption method. The symmetric encryption method is a secret-key algorithm, in which the sender and receiver must share the same key to perform encryption and decryption of the data. Whereas, asymmetric encryption method (also known as public-key cryptography) uses two keys, i.e., public and private key, associated with an entity which require to authenticate its identity electronically or encrypt data. The public key of each entity is published whereas, the corresponding private key is always kept secret to perform encryption or decryption of data. In literature, Ron Rivest, Adi Shamir and Leonard Adleman (RSA) [122], Data Encryption

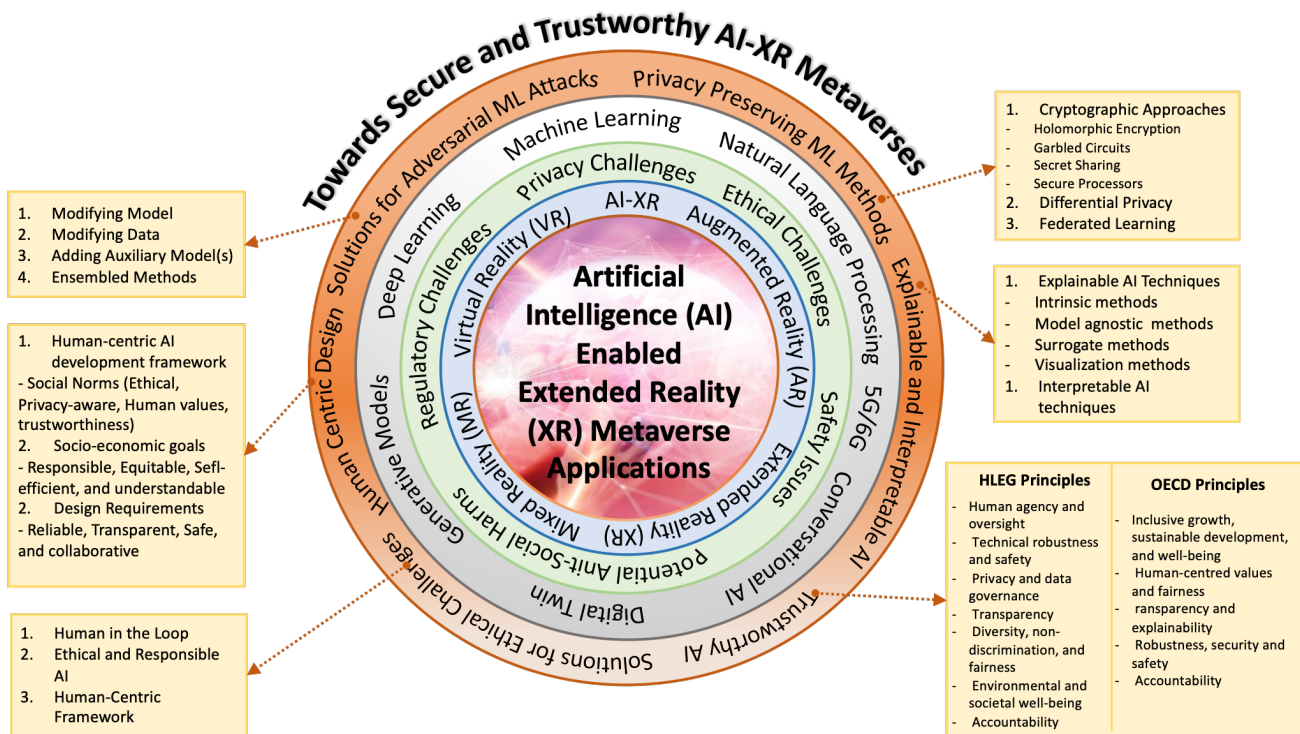


Fig. 11. An abstraction of different ML-associated challenges along with a taxonomy of various solutions that can be used to address those challenges.

Standard (DES) [123], Advanced Encryption Standard (AES) [124], and Secure Hash Algorithm (SHA) [125] are a few commonly used algorithms used for data encryption. Different cryptographic techniques can be employed to convert readable information to an encrypted state, which can be later used at the receiver end after performing decryption. Below we discuss some of the most commonly used encryption methods that can be used for the development of privacy-aware AI models.

*a) Homomorphic Encryption:* Homomorphic encryption (HE) is a computational approach that performs encryption while allowing computational tasks to be executed over encrypted data at the same time to ensure the privacy of the data. HE is defined as a public key cryptographic technique in which a pair of public and private keys is created to perform encryption and decryption operations on the data. The public key is used to encrypt the data, before sharing it with the third party for further computational tasks including training, and/or inference. Due to the homomorphic characteristics of this approach, the results can be decoded using the private key to visualize the results without showing them to third-party servers or unauthenticated users. In the ML literature, HE has been used for protecting the privacy of the users' data for different applications such as genome imputation [126], misinformation detection in text messages [127], etc. Specifically, the AI models are trained and inferred using encrypted training and testing data thus preserving the privacy of the sensitive data.

*b) Secure Multi-party Computation:* Secure multi-party computation (also known as secure computation) is a type of cryptographic technique that is focused on the development of collaborative methods to perform joint computation and

calculate a function over joint inputs while possessing those inputs in an isolated fashion. Contrary to the traditional cryptographic methods, where cryptography ensures confidentiality and integrity of communication or storage, while the adversary is outside the system of users, this approach protects the users' privacy from each other while performing ML-based tasks including training and inference activities.

*c) Garbled Circuits:* The idea of garbled circuits was first proposed by Yao in 1986 to perform two-party computation [128]. Garbled circuits can be used in a scenario where multiple parties are interested in performing some computation without sharing their data. Let's assume two parties (e.g., Alice and Bob for the sake of simplicity) want to perform some computation using garbled circuits. Alice will send his input and function in the form of a garbled circuit and Bob will utilize his garbled input with the garbled circuit to get the result of the required function, once he obtains it from Alice in an oblivious fashion. In [129], garbled circuits along with HE have been used to develop privacy-aware ML models, where the authors trained three classification models namely decision tree, Naïve Bayes, and hyperplane decision classifier using the encrypted data.

*d) Secret Sharing:* In secret sharing, multiple parties collaborate in the computation by sharing their secrets among them while holding a "share" of the individual secrets. The secret can only be reconstructed by combing all the individual shares kept by participating parties, otherwise, it will be useless. In the literature, the secret sharing technique has been successfully used for training AI models in a privacy-preserving way. For instance, Bonawitz et al. [130] used the secret sharing technique to train an ML model by aggregat-



ing model updates from multiple parties in a privacy-aware way. In a similar study [131], authors used this technique for the development of a privacy-aware ML-based emotion recognition system leveraging client-server architecture. In their proposed framework, the secret sharing technique was used for the communication of audio-visual data from the client side to the server, where an ensemble model based on a sparse autoencoder and a CNN model was used for the feature extraction from the collected data. The SVM classifier was then trained using the extracted features for the emotion recognition task. A secret sharing-based parallelized variant of principal component analysis (PCA) for preserving data privacy is presented in [132].

*e) Secure Processors:* Secure processors were pioneered by rogue software to protect sensitive code from being accessed by malicious actors at higher privilege levels. Secure processors are being used in different processors now to perform privacy-preserving operations, e.g., the Intel SGX processor. In [133], SGX processors were used to developing a data oblivious system for different ML techniques that include SVM, decision tree, matrix factorization, and k-mean clustering. The primary goal was to facilitate collaboration between multiple data proprietors performing the ML task on an SGX-empowered data center.

*2) Differential Privacy:* The idea of differential privacy is based on introducing noise in the data to protect sensitive information while ensuring the usefulness of the data after noise addition [134]. Differential privacy is defined in terms of the task-specific concept of neighbor datasets and it provides strong guarantees in ensuring the privacy of the data during algorithmic analysis [135]. Numerous differential privacy-based methods have been presented in the literature, such as differentially-private stochastic gradient descent (DP-SGD) [136], private-aggregation of teacher ensembles (PATE) [137], exponential noise based differential privacy-preserving methods to ensure privacy on large-scale data. These methods demonstrated better applicability in ML-based applications in various domains including intelligent transportation services, smart/virtual personal assistants, and smart healthcare services.

*3) Federated Learning:* Federated learning (FL) refers to a distributed-ML paradigm that is capable of learning global ML models without directly accessing and/or exchanging data from edge devices. Intuitively, basic FL-based methods consist of a collaborative learning framework where each participant such as an edge device, network node, and local server can independently train a model using its local data. These edge devices then share their model parameters with a server, which then performs aggregation of the parameters after receiving parameter updates from each edge device. Finally, the server updates the parameters of the global model and shares the updated parameters with all participants. The iterative process of FL is continuous until the desired criteria as been fulfilled, e.g., validation accuracy/loss or the maximum number of communication rounds. In this way, a global model is trained without requiring the actual data from the FL participants. Subsequently, this sharing mechanism allows ML-based systems to learn from large-scale diverse data and develop a global model. Such methods can demonstrate

better applicability in terms of dealing with sensitive data in various human-centered applications such as AI-XR metaverse applications. Despite the success of FL in training an ML model with reliable performance while maintaining the privacy of the actual data, different attacks can be realized on the model being trained using the FL paradigm, e.g., backdoor attacks [138], label flipping attacks [139], free-riding attacks [140], and poisoning attacks [141], etc. Also, it has been demonstrated that sensitive information can be extracted from the shared parameters in FL settings [142].

## B. Solutions to Combat Adversarial ML Attacks in AI-XR

In the literature, adversarially robust ML models have been mainly categorized into three categories [30]: (1) Data Modification; (2) Model Modification; and (3) Using Auxiliary Model. Moreover, a few methods leverage a hybrid approach in which multiple defensive techniques are used to develop adversarially robust ML models. Below we discuss the most prominent methods in each category and we refer the interested readers for more details about these methods to recent and comprehensive surveys that are specifically focused on adversarial ML [8], [45], [46], [143].

*1) Data modification:* Data modification methods work by modifying the input data during the training or inference phase to mitigate the effects of adversarial perturbation. A few famous data modification methods are briefly described below.

- *Adversarial Re-training:* This method was proposed by Goodfellow et al. [144] and it is considered to be a basic method for mitigating the effect of adversarial perturbation in the trained model. In this method, adversarial examples are augmented in the training data, which is then used to (re)-train the model. This method has been extensively used in the literature, however, a few research studies demonstrated that the models trained using this method are not robust against multiple attacks [145].
- *Feature Squeezing:* Xu et al. [146] presented a feature squeezing-based approach that aims to squeeze feature space of input that may be exploited in response to an adversary. In this regard, the heterogeneous feature vectors have been collectively joined into a single space to reduce available feature space. Although, the proposed defense method achieved significant performance against small perturbations. However, it was found less effective against iterative adversarial attacks [147].
- *Input Reconstruction:* Input reconstruction-based defense methods have been proposed to mitigate the effect of adversarial attacks. These methods transform adversarial examples into legitimate samples by cleaning adversarial noise using an appropriate technique, e.g., using an autoencoder to clean adversarial perturbations [148].

*2) Model modification:* Model modification methods aim at modifying the parameters of trained ML models to defuse the effect of adversarial attacks. The most commonly used model modification methods are described below.

- *Gradient Regularization:* This method allows complex neural networks to bring a partial surge in training

computational complexity to improve the performance of the network regardless of any prior knowledge about adversarial attacks. This idea was coined by Ross et al. [149] to improve the performance of CNN models on classification tasks. Though the proposed method achieved significant improvement in CNNs' robustness, it also increases the computational cost of models which prejudices the performance in real-world ML-based applications.

- *Defensive Distillation*: Distillation in a neural network was initially conceptualized by Hinton et al. [150] to establish knowledge sharing from a larger network to a smaller one. Later, Papernot et al. [151] extended this notion by developing a distillation-based defense mechanism against adversarial attacks, which is known as defensive distillation. In this method, the larger model is trained over hard labels to maximize accuracy while predicting the output probabilities of the baseline smaller model. This method is successful in mitigating the effect of small adversarial perturbation and it fails in the presence of strong adversarial perturbations, e.g., adversarial examples generated using C&W attack [152].
- *Network Verification*: In this method, certain properties of the ML/DL model are verified, e.g., validating the output of models, produced in response to the corresponding input samples. Katz et al. [153] presented ReLU and satisfiability modulo theory (SMT) based network verification method to make complex neural networks robust against adversarial examples. In a similar study, authors have proposed a scalable quantitative verification framework for DNNs to prove formal probabilistic property against adversarial attacks [154].

3) *Using Auxiliary Model*: Methods aiming to robustify ML models in this category use an additional model either for detection of adversarial examples or for clean adversarial perturbations. A few methods are described below.

- *Adversarial Detection*: In such methods, a detector model is used to differentiate between normal and adversarial inputs, e.g., a binary classifier [155].
- *Ensembling Defenses*: In this defense strategy, an ensemble of different defensive techniques is created to withstand different adversarial attacks. PixelDefend is the most famous ensemble defense method that consists of two defense approaches, i.e., input reconstruction and adversarial detection [156].
- *Using Generative Modeling*: These types of methods leverage different ML/DL-based generative models for cleaning adversarial noise in adversarial examples to project them back to the same data manifold.

### C. Solutions for AI-XR Transparency and Trust Challenges

The true potential of AI-based applications in AI-XR metaverse applications can only be realized when they are developed using fine-grained personal data for making personalized recommendations and predictions, which is only possible when users fully trust the underlying system. Therefore, addressing

the challenges related to the trustworthiness aspects of AI-XR metaverse applications is very important. From an AI perspective trustworthiness itself requires predictability, interpretability, explainability, safety, and robustness. Below we discuss different methods that can be used to accomplish trustworthiness in AI applications.

1) *Explainable and Interpretable AI*: An AI model is referred to as explainable if it can explain the ability of parameters to justify the results. Explainability makes the AI models transparent which ultimately helps in evaluating and understanding the results provided by the models. In recent years, substantial research efforts have been conducted to enhance explainability, trustworthiness, and interpretability in AI models. Fairness, Accountability, and Transparency in Machine Learning (FAT-ML) [157] and Defense Advanced Research Projects Agency (DARPA), explainable AI program [158] are the two famous research groups working in this context. The literature argues that explainable models can be the first step toward converting black-box AI models into white-box models [159].

Interpretable models refer to the models that explain themselves. In simple words, an AI model is said to be interpretable, if its decision against some input is logically understandable such as which factors influenced the AI model to reach that decision. In the literature, various methods have been presented to leverage interpretability in ML models. These methods ensure that the predictions of interpretable models are unbiased, which ultimately makes it easier to trust these systems in human society. It is worth noting that the terms interpretable and explainable are interchangeably used in the literature, however, they are different in terms of domain-specific definitions, moreover, there is no exact definition of these terms [160]. A detailed taxonomy of different explainable and interpretable AI methods can be found in [159], [160].

2) *Trustworthy AI*: The relevant literature emphasizes two famous sets of principles that can be used to attain trustworthy AI. One of them is developed by European Commission's AI High-Level Expert Group (HLEG) [161] and the other one is defined by Organisation for Economic Co-operation and Development (OECD) [162]. The following are the seven essential principles outlined in OECD: (1) Human agency and oversight; (2) Technical robustness and safety; (3) Privacy and data governance; (4) Transparency; (5) Diversity, non-discrimination, and fairness; (6) Environmental and societal well-being; and (7) Accountability.

Similarly, the following principles are outlined in HLEG: (1) Inclusive growth, sustainable development, and well-being; (2) Human-centred values and fairness; (3) Transparency and explainability; (4) Robustness, security, and safety; and (5) Accountability. One of the key noticeable insights from the above two principles set is that they mainly emphasized explainability, security, fairness, safety, and robustness aspects of AI. Therefore, these are the essential requirements that need to be fulfilled to develop trustworthy AI-based applications. In addition, we can see that these principles are essentially human-centric that respect ethical norms. As potential AI-XR metaverse applications will be more human-focused, therefore, the above-mentioned principles can be leveraged to develop

trustworthy AI-based applications for the metaverse.

#### D. Solutions for Ethical Challenges in AI-XR

1) *Human in the Loop*: The metaverse’s inherent complexity raises different security issues. For instance, it can be envisioned that the metaverse administrators will have to push automation, that is, to handle more tasks with algorithms, rather than with human operators, due to the requirement of managing a large number of users, applications, and services. The generated data will be much larger than those managed by the current Web platforms. Delegating tasks to algorithms, especially those implemented with state-of-the-art AI approaches is necessary to meet high-level efficiency and scalability. However, in the current version of social media and the Internet, we have even started to realize the implications of using algorithms for managing societally relevant tasks. Despite the significant performance, these algorithms suffer from various issues. Some authors writing on the governance of metaverse have proposed the use of a modular approach for the development of metaverse applications, as it allows adapting regulations to specific scenarios and then controlling the system accordingly [163].

2) *Ethical and Responsible AI*: To ensure socially desirable AI decisions, novel ways are required to be figured out to simultaneously minimize potential harms associated with the use of AI and its potential benefits. In this regard, the importance of taking an ethics-first approach towards the development of AI-based technologies becomes more plausible [164]. However, there are many challenges associated with the development of ethical AI pipelines due to distinct social norms and demographics of the human population, i.e., one ethical solution may be beneficial for a group of people but it is highly possible that it will not be suitable for another group on the same time. Therefore, customized solutions are required to address such issues that can consider the social norms of target users while making AI-based decisions. In this regard, different ethical guidelines can be leveraged that can be potentially used for the development of pro-social AI solutions. The literature shows a groundswell of interest in ensuring ethical and responsible AI [165].

#### E. Situational Awareness

Situational awareness can be defined as the capacity to understand information perceived from the surrounding environment. The literature argues that situational awareness is a crucial and effective tool for monitoring the security of complex systems like metaverse [166]. Situational awareness can be used at the local and global levels for threat monitoring in a single metaverse or across multiple metaverses, respectively. The feasibility and potential of this tool have been extensively studied in the literature focused on XR and VR technology. For instance, Woodward et al. [166] performed a literature review that focused on the design of information presentation in AR headsets to enhance users’ situational awareness. Authors in [167], performed immersive and realistic simulations to evaluate the effectiveness of audio-visual warning systems in increasing users’ situational awareness in accident situations

using VR. They demonstrated that VR can assist drivers to remain alert in emergency situations.

#### F. Human Centric Approach for AI-XR Development

Metaverse is essentially a human-centric application [168]. To realize the real social impact of different AI-XR metaverse applications, they should be analyzed and developed using human-centric design thinking. Metaverse service providers and developers must pay attention to key stakeholders (i.e., humans) by prioritizing and considering their social norms, i.e., dignity, justice, and rights, and supporting goals including creativity, self-efficacy, social connections, and responsibility. The aforementioned characteristics can be inherited in AI-XR metaverse applications by following three key concepts proposed in [169]. The first one is *Human-centric framework*—that guides the developers and researchers to ensure human-centric thinking about high-level two-dimensional control. Secondly, *Design metaphors*—which points out how two key goals of AI and social norms are both valuable. However, the stakeholders such as developers, researchers, policymakers, and business leaders must combine them both in developing metaverse applications to provide ultimate benefits to the users. Thirdly, *Governance Structures*—which ensures the bridge between the above-mentioned ethical principles and the practical measures needed to achieve the desired goals including reliable metaverse application development while ensuring cultural safety to increase privacy and trustworthiness of the users.

## VI. OPEN RESEARCH ISSUES

In this section, we highlight various open research issues that are particularly associated with the use of ML/DL models in different AI-XR metaverse applications.

#### A. Developing Generalizable Adversarial Defense Methods

Over the past few years, substantial research attention has been devoted to adversarial ML. However, the literature highlights that the attention devoted to developing adversarially robust ML/DL models is significantly less as compared to developing novel attack methodologies [115]. In the literature, different defensive techniques have been proposed to withstand adversarial ML attacks (as discussed above), however, each method only works in a specific setting and fails to withstand unseen and powerful attacks (consequently, fails to generalize across a wider class of attacks). On the other hand, the literature focused on adversarial ML shows that the diversity and severity of these attacks are increasing with each passing day. Therefore, the development of hybrid and universal defensive techniques is the need of the hour. In addition, it is required that the defense techniques should be developed while considering evolvable and adaptable adversaries (who can adapt their capabilities to break defense strategy). The threat of adversarial ML can be a major hurdle in the development of secure, safe, robust, and trustworthy AI-XR metaverse applications and if it remained unaddressed, can cause unintended severe consequences to users and society. It is highly recommended

to consider these aspects while developing ML/DL-empowered human-centric applications like the AI-XR metaverse. Moreover, the worst-case robustness test can be performed from an adversarial lens considering different attack surfaces in individual AI-XR metaverse application architecture.

### *B. Investigating Robustness of Privacy-Preserving Methods*

As discussed above, AI-XR metaverse applications will collect fine-grained data that may include personal attributes to provide personalized services (empowered by ML/DL models). The models trained with such data can be inferred to reconstruct privacy-related information that can be exploited to get intended outcomes and incentives. Although different privacy-preserving ML techniques have been proposed in the literature that has been shown quite successful in preserving data privacy, however, the literature demonstrates that meaningful information can still be inferred even if the presence of an appropriate privacy-preserving method. For example, it has been demonstrated that homomorphic encryption (one of the widely used encryption techniques) is vulnerable to model extraction attacks [170]. Similarly, Boenisch et al. [142] showed that sensitive information can be reconstructed from the shared parameters in FL. This suggests that the investigation of vulnerabilities and limitations of existing privacy-preserving methods can be a good step toward developing robust privacy-preserving methods. Ideally, it is required that the ML/DL models should be developed in such a way that they are by design privacy-aware, i.e., they should not be able to learn any privacy-related features from the data that could be compromised upon model inferences.

### *C. Developing Generalizable Explainable and Interpretable Techniques*

Another major limitation of DL models hindering their trustworthy applications in critical applications like the AI-XR metaverse is the lack of explainability and interpretability. This can also be exploited by adversarial agents to craft adversarial perturbations to realize attacks on different AI-XR metaverse applications. Although significant research interest has been devoted to the development of novel techniques to explain and interpret DL models, the literature shows that their application is limited to a certain data type or application [160]. While the AI-XR metaverse applications will have a complex architecture that will simultaneously use multi-modal data for making different intelligent decisions, existing explainability and interpretability techniques cannot be directly used for explaining and interpreting ML/DL-driven decisions. More work is required to create methods that can be generalized across different data types, models, and applications.

### *D. Developing Ethical Data Analysis Pipelines*

Current ML/DL models are not capable of considering different ethical norms that are necessary for human-centric applications like AI-XR metaverse applications. On contrary, these considerations are yet very important to maximize potential benefits and minimize associated harms to ensure safe,

robust, and fair data analysis. In AI-XR metaverse applications, different ML/DL models will be trained using massively large data collected by humans and their interactions with the real and virtual universe. While there is no guarantee that the AI decision will be ethically-committed because the data used for model training might contain data bias that will eventually result in biased decisions. Moreover, the outcomes of ML/DL models will be just a reflection of human behavior (including moral failures even if they are not intentionally committed). Therefore, to increase the trust of different stakeholders involved in AI-XR metaverse applications (particularly, end users) and to provide them with a sense of safety, fairness, and accountability, it is highly desirable to develop novel techniques to ensure fair and ethical data analysis empowered by various ML/DL techniques.

### *E. Pushing AI on Edge: Embedded ML*

One of the feasible approaches to preserve the privacy of AI-XR metaverse users will be to deploy ML/DL models on their smart devices, e.g., smartphones, AR/VR gadgets, tablets, etc. By doing so models can be developed and inferred on their devices without requiring to transmit data to a central cloud. We envision that various AI-XR metaverse applications will potentially adopt embedded ML or edge-enabled ML due to the proliferation of different enabling gadgets and smart devices. However, numerous challenges related to underlying hardware computing capabilities will arise when sufficiently large ML/DL models will be deployed on resource-constrained devices. Also, the literature argues that the research on enabling edge AI is at its early stages of development [171]. Therefore, it is worth investigating the feasibility and potential of deploying ML/DL models on embedded devices to ultimately develop secure, private, and robust systems to provide personalized services in AI-XR metaverse applications. We refer interested readers to a recent survey on analyzing the notion of edge-enabled metaverse applications for a more comprehensive discussion on the topic and various challenges [172].

### *F. AI-XR Metaverse Specific Security Solutions*

The future AI-XR metaverse will have a complex structure and will be a combination of various enabling (complex) technologies (that possess their associated challenges related to privacy and security, e.g., adversarial ML). Moreover, the massive connectivity of numerous entities (users, services providers, organizations, etc.) along with the decentralization will even worsen the enormity of security and privacy in AI-XR metaverse applications. Individual vulnerabilities associated with each technology can be exploited to realize a more powerful attack to halt or get control of some services or the entire metaverse. If such vulnerabilities are left unaddressed, they will eventually lead to novel challenges thus making it challenging to ensure the secure, safe, and robust operation of metaverse services. Therefore, it is very crucial to understand such challenges and develop customized defense solutions to protect AI-XR metaverse applications and services in general.

## VII. CONCLUSIONS

In this paper, we have analyzed various security, privacy, and trustworthiness challenges associated with the use of different machine learning (ML) and deep learning (DL) techniques in artificial intelligence and extended reality (AI-XR) metaverse applications. Specifically, considering the layered architecture of the metaverse, we developed a pipeline and highlighted different potential ML/DL use cases along with identifying various vulnerabilities associated with their application. Furthermore, we provide a comprehensive overview of these challenges and discuss potential solutions that could be used to overcome such issues. To accentuate the implications of adversarial threats, we designed a customized case study (considering a prospective AI-XR metaverse application) and analyzed its security and privacy aspects. Finally, we discussed various open research issues that require further investigation. We envision that our work on this crucial topic will provide a one-stop solution to interested researchers who aim to develop secure, robust, and trustworthy AI-XR applications.

## REFERENCES

- [1] R. Cheng, N. Wu, S. Chen, and B. Han, "Will metaverse be NextG internet? vision, hype, and reality," *arXiv preprint arXiv:2201.12894*, 2022.
- [2] B. Falchuk, S. Loeb, and R. Neff, "The social metaverse: Battle for privacy," *IEEE Technology and Society Magazine*, vol. 37, no. 2, pp. 52–61, 2018.
- [3] J. A. De Guzman, K. Thilakarathna, and A. Seneviratne, "Security and privacy approaches in mixed reality: A literature survey," *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–37, 2019.
- [4] H. Ning, H. Wang, Y. Lin, W. Wang, S. Dhelim, F. Farha, J. Ding, and M. Daneshmand, "A survey on metaverse: the state-of-the-art, technologies, applications, and challenges," *arXiv preprint arXiv:2111.09673*, 2021.
- [5] R. Di Pietro and S. Cresci, "Metaverse: Security and privacy issues," in *2021 Third IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*. IEEE, 2021, pp. 281–288.
- [6] T. Huynh-The, Q.-V. Pham, X.-Q. Pham, T. T. Nguyen, Z. Han, and D.-S. Kim, "Artificial intelligence for the metaverse: A survey," *arXiv preprint arXiv:2202.10336*, 2022.
- [7] R. Zhao, Y. Zhang, Y. Zhu, R. Lan, and Z. Hua, "Metaverse: Security and privacy concerns," *arXiv preprint arXiv:2203.03854*, 2022.
- [8] Y. Wang, Z. Su, N. Zhang, R. Xing, D. Liu, T. H. Luan, and X. Shen, "A survey on metaverse: Fundamentals, security, and privacy," *IEEE Communications Surveys & Tutorials*, 2022.
- [9] O. Halabi, S. Balakrishnan, S. P. Dakua, N. Navab, and M. Warfa, "Virtual and Augmented Reality in Surgery," in *The Disruptive Fourth Industrial Revolution*. Springer, 2020, no. July, pp. 257–285.
- [10] S. B. Far and A. I. Rad, "Applying digital twins in metaverse: User interface, security and privacy challenges," *Journal of Metaverse*, vol. 2, no. 1, pp. 8–16, 2022.
- [11] D. Reiners, M. R. Davahli, W. Karwowski, and C. Cruz-Neira, "The combination of artificial intelligence and extended reality: A systematic review," *Frontiers in Virtual Reality*, p. 114, 2021.
- [12] D. W. Otter, J. R. Medina, and J. K. Kalita, "A survey of the usages of deep learning for natural language processing," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 2, pp. 604–624, 2020.
- [13] M. Chen, U. Challita, W. Saad, C. Yin, and M. Debbah, "Artificial neural networks-based machine learning for wireless networks: A tutorial," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3039–3071, 2019.
- [14] C. She, R. Dong, Z. Gu, Z. Hou, Y. Li, W. Hardjawana, C. Yang, L. Song, and B. Vucetic, "Deep learning for ultra-reliable and low-latency communications in 6g networks," *IEEE network*, vol. 34, no. 5, pp. 219–225, 2020.
- [15] M. Alsenwi, N. H. Tran, M. Bennis, S. R. Pandey, A. K. Bairagi, and C. S. Hong, "Intelligent resource slicing for embb and urllc coexistence in 5g and beyond: A deep reinforcement learning based approach," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, pp. 4585–4600, 2021.
- [16] C. Luo, J. Ji, Q. Wang, X. Chen, and P. Li, "Channel state information prediction for 5g wireless communications: A deep learning approach," *IEEE Transactions on Network Science and Engineering*, vol. 7, no. 1, pp. 227–236, 2018.
- [17] G. B. Tunze, T. Huynh-The, J.-M. Lee, and D.-S. Kim, "Sparsely connected cnn for efficient automatic modulation recognition," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 12, pp. 15557–15568, 2020.
- [18] A. Cannavo and F. Lamberti, "How blockchain, virtual reality, and augmented reality are converging, and why," *IEEE Consumer Electronics Magazine*, vol. 10, no. 5, pp. 6–13, 2020.
- [19] Q. Yang, Y. Zhao, H. Huang, Z. Xiong, J. Kang, and Z. Zheng, "Fusing blockchain and ai with metaverse: A survey," *IEEE Open Journal of the Computer Society*, vol. 3, pp. 122–136, 2022.
- [20] S. Tanwar, Q. Bhatia, P. Patel, A. Kumari, P. K. Singh, and W.-C. Hong, "Machine learning adoption in blockchain-based smart applications: The challenges, and a way forward," *IEEE Access*, vol. 8, pp. 474–488, 2019.
- [21] F. Tao, H. Zhang, A. Liu, and A. Y. Nee, "Digital twin in industry: State-of-the-art," *IEEE Transactions on industrial informatics*, vol. 15, no. 4, pp. 2405–2415, 2018.
- [22] D. Chen, D. Wang, Y. Zhu, and Z. Han, "Digital twin for federated analytics using a bayesian approach," *IEEE Internet of Things Journal*, vol. 8, no. 22, pp. 16301–16312, 2021.
- [23] L.-H. Lee, T. Braud, P. Zhou, L. Wang, D. Xu, Z. Lin, A. Kumar, C. Bermejo, and P. Hui, "All one needs to know about metaverse: A complete survey on technological singularity, virtual ecosystem, and research agenda," *arXiv preprint arXiv:2110.05352*, 2021.
- [24] S. Wasserkrug, A. Gal, and O. Etzion, "Inference of security hazards from event composition based on incomplete or uncertain information," *IEEE transactions on knowledge and data engineering*, vol. 20, no. 8, pp. 1111–1114, 2008.
- [25] J. Wei, J. Li, Y. Lin, and J. Zhang, "Ldp-based social content protection for trending topic recommendation," *IEEE Internet of Things Journal*, vol. 8, no. 6, pp. 4353–4372, 2020.
- [26] J. Shang, S. Chen, J. Wu, and S. Yin, "Arspy: Breaking location-based multi-player augmented reality application for user location tracking," *IEEE Transactions on Mobile Computing*, 2020.
- [27] J. D. Lee and K. A. See, "Trust in automation: Designing for appropriate reliance," *Human factors*, vol. 46, no. 1, pp. 50–80, 2004.
- [28] H. Ali, M. S. Khan, A. Al-Fuqaha, and J. Qadir, "Tamp-X: Attacking explainable natural language classifiers through tampered activations," *Computers & Security*, p. 102791, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167404822001857>
- [29] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.
- [30] A. Qayyum, J. Qadir, M. Bilal, and A. Al-Fuqaha, "Secure and robust machine learning for healthcare: A survey," *IEEE Reviews in Biomedical Engineering*, vol. 14, pp. 156–180, 2020.
- [31] S. Latif, A. Qayyum, M. Usama, J. Qadir, A. Zwitter, and M. Shahzad, "Caveat emptor: the risks of using big data for human development," *IEEE technology and society magazine*, vol. 38, no. 3, pp. 82–90, 2019.
- [32] P. Kürtünlüoğlu, B. Akdik, and E. Karaarslan, "Security of virtual reality authentication methods in metaverse: An overview," *arXiv preprint arXiv:2209.06447*, 2022.
- [33] E. Ntoutsi *et al.*, "Bias in data-driven artificial intelligence systems—an introductory survey," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 10, no. 3, p. e1356, 2020.
- [34] A. Olteanu, C. Castillo, F. Diaz, and E. Kiciman, "Social data: Biases, methodological pitfalls, and ethical boundaries," *Frontiers in Big Data*, vol. 2, p. 13, 2019.
- [35] H. Suresh and J. V. Gutttag, "A framework for understanding unintended consequences of machine learning," *arXiv preprint arXiv:1901.10002*, 2019.
- [36] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," *arXiv preprint arXiv:1908.09635*, 2019.
- [37] S. Vallor, *Technology and the virtues: A philosophical guide to a future worth wanting*. Oxford University Press, 2016.
- [38] J. Lanier, *Ten arguments for deleting your social media accounts right now*. Random House, 2018.

- [39] Marr, *Extended Reality in Practice*. Wiley, 2021.
- [40] U. H. of Commons DCMS Committee *et al.*, “Immersive and addictive technologies. parliament. uk,” 2019.
- [41] M. Flintham, C. Karner, K. Bachour, H. Creswick, N. Gupta, and S. Moran, “Falling for fake news: investigating the consumption of news via social media,” in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, 2018, p. 376.
- [42] L. Buck and R. McDonnell, “Security and privacy in the metaverse: The threat of the digital human,” *Proceedings of the 1st Workshop on Novel Challenges of Safety, Security and Privacy in Extended Reality*, 2022.
- [43] K. Shahriari and M. Shahriari, “Ieee standard review—ethically aligned design: A vision for prioritizing human wellbeing with artificial intelligence and autonomous systems,” in *2017 IEEE Canada International Humanitarian Technology Conference (IHTC)*. IEEE, 2017, pp. 197–201.
- [44] L. B. Rosenberg, “Regulation of the metaverse: A roadmap,” in *6th International Conference on Virtual and Augmented Reality Simulations (ICVARs 2022)*, 2022.
- [45] A. Qayyum, M. Usama, J. Qadir, and A. Al-Fuqaha, “Securing connected & autonomous vehicles: Challenges posed by adversarial machine learning and the way forward,” *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 998–1026, 2020.
- [46] N. Akhtar, A. Mian, N. Kardan, and M. Shah, “Advances in adversarial attacks and defenses in computer vision: A survey,” *IEEE Access*, vol. 9, pp. 155 161–155 196, 2021.
- [47] B. Biggio, B. Nelson, and P. Laskov, “Poisoning attacks against support vector machines,” in *29th International Conference on Machine Learning*. ArXiv e-prints, 2012, pp. 1807–1814.
- [48] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto, and F. Roli, “Evasion attacks against machine learning at test time,” in *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 2013, pp. 387–402.
- [49] F. Khalid, H. Ali, M. A. Hanif, S. Rehman, R. Ahmed, and M. Shafique, “Fadec: A fast decision-based attack for adversarial machine learning,” in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–8.
- [50] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, “Badnets: Evaluating backdooring attacks on deep neural networks,” *IEEE Access*, vol. 7, pp. 47 230–47 244, 2019.
- [51] H. Ali, S. Nepal, S. S. Kanhere, and S. Jha, “Has-nets: A heal and select mechanism to defend dnns against backdoor attacks for data collection scenarios,” *arXiv preprint arXiv:2012.07474*, 2020.
- [52] H. Ali, M. S. Khan, A. AlGhadhban, M. Alazmi, A. Alzamil, K. Al-Utaibi, and J. Qadir, “All your fake detector are belong to us: Evaluating adversarial robustness of fake-news detectors under black-box settings,” *IEEE Access*, vol. 9, pp. 81 678–81 692, 2021.
- [53] H. Ali, M. S. Khan, A. AlGhadhban, M. Alazmi, A. Alzamil, K. Al-Utaibi, and J. Qadir, “Con-detect: Detecting adversarially perturbed natural language inputs to deep classifiers through holistic analysis,” *TechRxiv*, 2022.
- [54] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [55] M. A. Butt, A. M. Khattak, S. Shafique, B. Hayat, S. Abid, K.-I. Kim, M. W. Ayub, A. Sajid, and A. Adnan, “Convolutional neural network based vehicle classification in adverse illuminous conditions for intelligent transportation systems,” *Complexity*, vol. 2021, 2021.
- [56] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.
- [57] M. A. Butt and F. Riaz, “Carl-d: a vision benchmark suite and large scale dataset for vehicle detection and scene segmentation,” *Signal Processing: Image Communication*, vol. 104, p. 116667, 2022.
- [58] W. Xia, Y. Yang, J.-H. Xue, and B. Wu, “Tedigan: Text-guided diverse face image generation and manipulation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2256–2265.
- [59] F. Khalid, H. Ali, H. Tariq, M. A. Hanif, S. Rehman, R. Ahmed, and M. Shafique, “Qusecnets: Quantization-based defense mechanism for securing deep neural network against adversarial attacks,” in *2019 IEEE 25th International Symposium on On-Line Testing and Robust System Design (IOLTS)*. IEEE, 2019, pp. 182–187.
- [60] H. Ali, F. Khalid, H. A. Tariq, M. A. Hanif, R. Ahmed, and S. Rehman, “Sscnets: Robustifying dnns using secure selective convolutional filters,” *IEEE Design & Test*, vol. 37, no. 2, pp. 58–65, 2019.
- [61] M. Usama, J. Qadir, and A. Al-Fuqaha, “Adversarial attacks on cognitive self-organizing networks: The challenge and the way forward,” in *2018 IEEE 43rd Conference on Local Computer Networks Workshops (LCN Workshops)*. IEEE, 2018, pp. 90–97.
- [62] M. Usama, J. Qadir, M. A. Imran *et al.*, “Adversarial ml attack on self organizing cellular networks,” in *2019 UK/China Emerging Technologies (UCET)*. IEEE, 2019, pp. 1–5.
- [63] G. Goswami, N. Ratha, A. Agarwal, R. Singh, and M. Vatsa, “Unravelling robustness of deep learning based face recognition against adversarial attacks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [64] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, “Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition,” in *Proceedings of the 2016 acm sigsac conference on computer and communications security*, 2016, pp. 1528–1540.
- [65] M. Shen, Z. Liao, L. Zhu, K. Xu, and X. Du, “Vla: A practical visible light-based attack on face recognition systems in physical world,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 3, pp. 1–19, 2019.
- [66] E. Chatzikyriakidis, C. Papaioannidis, and I. Pitas, “Adversarial face de-identification,” in *2019 IEEE International conference on image processing (ICIP)*. IEEE, 2019, pp. 684–688.
- [67] A. Dabouei, S. Soleymani, J. Dawson, and N. Nasrabadi, “Fast geometrically-perturbed adversarial faces,” in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1979–1988.
- [68] Y. Zhong and W. Deng, “Towards transferable adversarial attack against deep face recognition,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1452–1466, 2020.
- [69] Y. Dong, H. Su, B. Wu, Z. Li, W. Liu, T. Zhang, and J. Zhu, “Efficient decision-based black-box adversarial attacks on face recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7714–7722.
- [70] E. Wenger, J. Passananti, A. N. Bhagoji, Y. Yao, H. Zheng, and B. Y. Zhao, “Backdoor attacks against deep learning systems in the physical world,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6206–6215.
- [71] M. Xue, C. He, J. Wang, and W. Liu, “Backdoors hidden in facial features: a novel invisible backdoor attack against face recognition systems,” *Peer-to-Peer Networking and Applications*, vol. 14, no. 3, pp. 1458–1474, 2021.
- [72] H. Zhang, W. Zhou, and H. Li, “Contextual adversarial attacks for object detection,” in *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2020, pp. 1–6.
- [73] M. Lee and Z. Kolter, “On physical adversarial patches for object detection,” *arXiv preprint arXiv:1906.11897*, 2019.
- [74] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille, “Adversarial examples for semantic segmentation and object detection,” in *Proceedings of the IEEE ICCV*, 2017, pp. 1369–1378.
- [75] X. Wei, S. Liang, N. Chen, and X. Cao, “Transferable adversarial attacks for image and video object detection,” *arXiv preprint arXiv:1811.12641*, 2018.
- [76] Y. Wang, Y.-a. Tan, W. Zhang, Y. Zhao, and X. Kuang, “An adversarial attack on dnn-based black-box object detectors,” *Journal of Network and Computer Applications*, vol. 161, p. 102634, 2020.
- [77] T. Wu, T. Wang, V. Schwag, S. Mahloujifar, and P. Mittal, “Just rotate it: Deploying backdoor attacks via rotation transformation,” *arXiv preprint arXiv:2207.10825*, 2022.
- [78] X. Wang, M. Cai, F. Sohel, N. Sang, and Z. Chang, “Adversarial point cloud perturbations against 3d object detection in autonomous driving systems,” *Neurocomputing*, vol. 466, pp. 27–36, 2021.
- [79] C. Xiang, C. R. Qi, and B. Li, “Generating 3d adversarial point clouds,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9136–9144.
- [80] A. Hamdi, S. Rojas, A. Thabet, and B. Ghanem, “Advpc: Transferable adversarial perturbations on 3d point clouds,” in *European Conference on Computer Vision*. Springer, 2020, pp. 241–257.
- [81] E. Meloni, M. Tiezzi, L. Pasqualini, M. Gori, and S. Melacci, “Messing up 3d virtual environments: Transferable adversarial 3d objects,” in *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2021, pp. 1–8.
- [82] X. Li, Z. Chen, Y. Zhao, Z. Tong, Y. Zhao, A. Lim, and J. T. Zhou, “Pointba: Towards backdoor attacks in 3d point cloud,” in *Proceedings of the IEEE ICCV*, 2021, pp. 16 492–16 501.

- [83] A. Arnab, O. Miksik, and P. H. Torr, "On the robustness of semantic segmentation models to adversarial attacks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2018, pp. 888–897.
- [84] J. Hendrik Metzen, M. Chaithanya Kumar, T. Brox, and V. Fischer, "Universal adversarial perturbations against semantic image segmentation," in *Proceedings of the IEEE ICCV*. IEEE, 2017, pp. 2755–2764.
- [85] Y. Li, Y. Li, Y. Lv, Y. Jiang, and S.-T. Xia, "Hidden backdoor attack against semantic segmentation models," *arXiv preprint arXiv:2103.04038*, 2021.
- [86] Y. Feng, B. Ma, J. Zhang, S. Zhao, Y. Xia, and D. Tao, "Fiba: Frequency-injection based backdoor attack in medical image analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 20 876–20 885.
- [87] L. Muñoz-González, B. Biggio, A. Demontis, A. Paudice, V. Wongrasamee, E. C. Lupu, and F. Roli, "Towards poisoning of deep learning algorithms with back-gradient optimization," in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 27–38.
- [88] J. Steinhardt, P. W. W. Koh, and P. S. Liang, "Certified defenses for data poisoning attacks," *Advances in neural information processing systems*, vol. 30, 2017.
- [89] A. Shafahi, W. R. Huang, M. Najibi, O. Suci, C. Studer, T. Dumitras, and T. Goldstein, "Poison frogs! targeted clean-label poisoning attacks on neural networks," *Advances in neural information processing systems*, vol. 31, 2018.
- [90] X. Zhang, J. Zhang, Z. Chen, and K. He, "Crafting adversarial examples for neural machine translation," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2021, pp. 1967–1977.
- [91] N. Boucher, I. Shumailov, R. Anderson, and N. Papernot, "Bad characters: Imperceptible nlp attacks," in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 1987–2004.
- [92] L. Li, R. Ma, Q. Guo, X. Xue, and X. Qiu, "Bert-attack: Adversarial attack against bert using bert," *arXiv preprint arXiv:2004.09984*, 2020.
- [93] D. Jin, Z. Jin, J. T. Zhou, and P. Szolovits, "Is bert really robust? a strong baseline for natural language attack on text classification and entailment," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 8018–8025.
- [94] R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, and Y. Choi, "Defending against neural fake news," *Advances in neural information processing systems*, vol. 32, 2019.
- [95] X. Pan, M. Zhang, B. Sheng, J. Zhu, and M. Yang, "Hidden trigger backdoor attack on {NLP} models via linguistic style manipulation," in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 3611–3628.
- [96] S. Garg and G. Ramakrishnan, "Bae: Bert-based adversarial examples for text classification," *arXiv preprint arXiv:2004.01970*, 2020.
- [97] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "Hotflip: White-box adversarial examples for text classification," *arXiv preprint arXiv:1712.06751*, 2017.
- [98] J. Li, S. Ji, T. Du, B. Li, and T. Wang, "Textbugger: Generating adversarial text against real-world applications," *arXiv preprint arXiv:1812.05271*, 2018.
- [99] X. Chen, A. Salem, M. Backes, S. Ma, and Y. Zhang, "Badnlp: Backdoor attacks against nlp models," in *ICML 2021 Workshop on Adversarial Machine Learning*, 2021.
- [100] S. Irtiza, L. Khan, and K. W. Hamlen, "Sentmod: Hidden backdoor attack on unstructured textual data," in *2022 IEEE 8th Intl Conference on Big Data Security on Cloud (BigDataSecurity), IEEE Intl Conference on High Performance and Smart Computing (HPSC) and IEEE Intl Conference on Intelligent Data and Security (IDS)*. IEEE, 2022, pp. 224–231.
- [101] M. Usama, M. Asim, S. Latif, J. Qadir *et al.*, "Generative adversarial networks for launching and thwarting adversarial attacks on network intrusion detection systems," in *2019 15th international wireless communications & mobile computing conference (IWCMC)*. IEEE, 2019, pp. 78–83.
- [102] M. Usama, J. Qadir, and A. Al-Fuqaha, "Black-box adversarial ml attack on modulation classification," *arXiv preprint arXiv:1908.00635*, 2019.
- [103] M. Usama, I. Ilahi, J. Qadir, R. N. Mitra, and M. K. Marina, "Examining machine learning for 5g and beyond through an adversarial lens," *IEEE Internet Computing*, vol. 25, no. 2, pp. 26–34, 2021.
- [104] O. Ibitoye, R. Abou-Khamis, A. Matrawy, and M. O. Shafiq, "The threat of adversarial attacks on machine learning in network security—a survey," *arXiv preprint arXiv:1911.02621*, 2019.
- [105] J. Aiken and S. Scott-Hayward, "Investigating adversarial attacks against network intrusion detection systems in sdns," in *2019 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*. IEEE, 2019, pp. 1–7.
- [106] M. Sadeghi and E. G. Larsson, "Adversarial attacks on deep-learning based radio signal classification," *IEEE Wireless Communications Letters*, vol. 8, no. 1, pp. 213–216, 2018.
- [107] Y. E. Sagduyu, Y. Shi, and T. Erpek, "Iot network security from the perspective of adversarial deep learning," in *2019 16th Annual International Conference on Sensing, Communication, and Networking*. IEEE, 2019, pp. 1–9.
- [108] Y. Wang, Z. Su, J. Ni, N. Zhang, and X. Shen, "Blockchain-empowered space-air-ground integrated networks: Opportunities, challenges, and solutions," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 1, pp. 160–209, 2021.
- [109] C. T. Nguyen, D. T. Hoang, D. N. Nguyen, and E. Dutkiewicz, "Metachain: A novel blockchain-based framework for metaverse applications," *arXiv preprint arXiv:2201.00759*, 2021.
- [110] E. Bertino and N. Islam, "Botnets and internet of things security," *Computer*, vol. 50, no. 2, pp. 76–79, 2017.
- [111] Y. Chen, X. Gong, Q. Wang, X. Di, and H. Huang, "Backdoor attacks and defenses for deep neural networks in outsourced cloud environments," *IEEE Network*, 2020.
- [112] T. S. Sethi and M. Kantardzic, "Data driven exploratory attacks on black box classifiers in adversarial domains," *Neurocomputing*, vol. 289, pp. 129–143, 2018.
- [113] Z. Yang, J. Zhang, E.-C. Chang, and Z. Liang, "Neural network inversion in adversarial setting via background knowledge alignment," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 225–240.
- [114] M. Kesarwani, B. Mukhoty, V. Arya, and S. Mehta, "Model extraction warning in MLaaS paradigm," in *Proceedings of the 34th Annual Computer Security Applications Conference*. ACM, 2018, pp. 371–380.
- [115] A. Qayyum, A. Ijaz, M. Usama, W. Iqbal, J. Qadir, Y. Elkhatib, and A. Al-Fuqaha, "Securing machine learning in the cloud: A systematic review of cloud machine learning security," *Frontiers in big Data*, vol. 3, p. 587139, 2020.
- [116] Y. Qin, N. Carlini, G. Cottrell, I. Goodfellow, and C. Raffel, "Imperceptible, robust, and targeted adversarial examples for automatic speech recognition," in *International conference on machine learning*. PMLR, 2019, pp. 5231–5240.
- [117] W. Xu, Y. Qi, and D. Evans, "Automatically evading classifiers," in *Proceedings of the 2016 network and distributed systems symposium*, vol. 10, 2016.
- [118] V. Fischer, M. C. Kumar, J. H. Metzen, and T. Brox, "Adversarial examples for semantic image segmentation," *arXiv preprint arXiv:1703.01101*, 2017.
- [119] Y. Xie, C. Shi, Z. Li, J. Liu, Y. Chen, and B. Yuan, "Real-time, universal, and robust adversarial attacks against speaker recognition systems," in *ICASSP 2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2020, pp. 1738–1742.
- [120] P. Casey, I. Baggili, and A. Yarramreddy, "Immersive virtual reality attacks and the human joystick," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 2, pp. 550–562, 2019.
- [121] S. Valluripally, A. Gulhane, R. Mitra, K. A. Hoque, and P. Calyam, "Attack trees for security and privacy in social virtual reality learning environments," in *2020 IEEE 17th Annual Consumer Communications & Networking Conference (CCNC)*. IEEE, 2020, pp. 1–9.
- [122] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Communications of the ACM*, vol. 21, no. 2, pp. 120–126, 1978.
- [123] D. E. Standard *et al.*, "Data encryption standard," *Federal Information Processing Standards Publication*, vol. 112, 1999.
- [124] N. F. Pub, "197: Advanced encryption standard (AES)," *Federal information processing standards publication*, vol. 197, no. 441, p. 0311, 2001.
- [125] J. H. Burrows, "Secure hash standard," Department of Commerce Washington DC, Tech. Rep., 1995.
- [126] E. Sarkar, E. Chielle, G. Gürsoy, O. Mazonka, M. Gerstein, and M. Maniatakos, "Fast and scalable private genotype imputation using machine learning and partially homomorphic encryption," *IEEE Access*, vol. 9, pp. 93 097–93 110, 2021.
- [127] H. Ali, R. T. Javed, A. Qayyum, A. AlGhadhban, M. Alazmi, A. Alzamil, K. Al-utaibi, and J. Qadir, "Spam-das: Secure and privacy-aware misinformation detection as a service," *TechRxiv*, 2022.

- [128] A. C.-C. Yao, "How to generate and exchange secrets," in *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*. IEEE, 1986, pp. 162–167.
- [129] R. Bost, R. A. Popa, S. Tu, and S. Goldwasser, "Machine learning classification over encrypted data," in *NDSS*, vol. 4324, 2015, p. 4325.
- [130] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2017, pp. 1175–1191.
- [131] M. S. Hossain and G. Muhammad, "Emotion recognition using secure edge and cloud computing," *Information Sciences*, vol. 504, pp. 589–601, 2019.
- [132] D. Bogdanov, L. Kamm, S. Laur, and V. Sokk, "Implementation and evaluation of an algorithm for cryptographically private principal component analysis on genomic data," *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 15, no. 5, pp. 1427–1432, 2018.
- [133] O. Ohrimenko, F. Schuster, C. Fournet, A. Mehta, S. Nowozin, K. Vaswani, and M. Costa, "Oblivious multi-party machine learning on trusted processors," in *25th USENIX Security Symposium (USENIX Security 16)*, 2016, pp. 619–636.
- [134] C. Dwork, "Differential privacy," *Encyclopedia of Cryptography and Security*, pp. 338–340, 2011.
- [135] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016, pp. 308–318.
- [136] J. Du, S. Li, M. Feng, and S. Chen, "Dynamic differential-privacy preserving SGD," *arXiv preprint arXiv:2111.00173*, 2021.
- [137] Q. Zhang, J. Ma, J. Lou, L. Xiong, and X. Jiang, "Towards training robust private aggregation of teacher ensembles under noisy labels," in *2020 IEEE international conference on big data (big data)*. IEEE, 2020, pp. 1103–1110.
- [138] C. Xie, K. Huang, P.-Y. Chen, and B. Li, "Dba: Distributed backdoor attacks against federated learning," in *International Conference on Learning Representations*, 2019.
- [139] A. Qayyum, M. U. Janjua, and J. Qadir, "Making federated learning robust to adversarial attacks by learning data and model association," *Computers & Security*, vol. 121, p. 102827, 2022.
- [140] J. Lin, M. Du, and J. Liu, "Free-riders in federated learning: Attacks and defenses," *arXiv preprint arXiv:1911.12560*, 2019.
- [141] A. Ali, I. Ilahi, A. Qayyum, I. Mohammed, A. Al-Fuqaha, and J. Qadir, "Incentive-driven federated learning and associated security challenges: A systematic review," *TechRxiv*, 2021.
- [142] F. Boenisch, A. Dziejzic, R. Schuster, A. S. Shamsabadi, I. Shumailov, and N. Papernot, "When the curious abandon honesty: Federated learning is not private," *arXiv preprint arXiv:2112.02918*, 2021.
- [143] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE transactions on neural networks and learning systems*, 2019.
- [144] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [145] F. Tramer and D. Boneh, "Adversarial training and robustness for multiple perturbations," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [146] W. Xu, D. Evans, and Y. Qi, "Feature squeezing: Detecting adversarial examples in deep neural networks," *arXiv preprint arXiv:1704.01155*, 2017.
- [147] W. He, J. Wei, X. Chen, N. Carlini, and D. Song, "Adversarial example defense: Ensembles of weak defenses are not strong," in *11th {USENIX} workshop on offensive technologies ({WOOT} 17)*, 2017.
- [148] S. Gu and L. Rigazio, "Towards deep neural network architectures robust to adversarial examples," *arXiv preprint arXiv:1412.5068*, 2014.
- [149] A. S. Ross and F. Doshi-Velez, "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [150] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [151] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *2016 IEEE symposium on security and privacy (SP)*. IEEE, 2016, pp. 582–597.
- [152] N. Carlini and D. Wagner, "Adversarial examples are not easily detected: Bypassing ten detection methods," in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 3–14.
- [153] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer, "Reluplex: An efficient smt solver for verifying deep neural networks," in *International Conference on Computer Aided Verification*. Springer, 2017, pp. 97–117.
- [154] T. Baluta, Z. L. Chua, K. S. Meel, and P. Saxena, "Scalable quantitative verification for deep neural networks," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, 2021, pp. 312–323.
- [155] J. Lu, T. Issararon, and D. Forsyth, "Safetynet: Detecting and rejecting adversarial examples robustly," in *Proceedings of the IEEE ICCV*, 2017, pp. 446–454.
- [156] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman, "Pixeldefend: Leveraging generative models to understand and defend against adversarial examples," in *International Conference on Learning Representations (ICLR)*, 2018. [Online]. Available: <https://openreview.net/forum?id=rJUYGxbCW>
- [157] M. FAT, "Fairness, accountability, and transparency in machine learning," *Retrieved December*, vol. 24, p. 2018, 2018.
- [158] D. Gunning, "Explainable artificial intelligence (XAI)," *Defense Advanced Research Projects Agency (DARPA), nd Web*, vol. 2, no. 2, 2017.
- [159] A. Adadi and M. Berrada, "Peeking inside the black-box: a survey on explainable artificial intelligence (xai)," *IEEE access*, vol. 6, pp. 52 138–52 160, 2018.
- [160] K. Rasheed, A. Qayyum, M. Ghaly, A. Al-Fuqaha, A. Razi, and J. Qadir, "Explainable, trustworthy, and ethical machine learning for healthcare: A survey," *Computers in Biology and Medicine*, p. 106043, 2022.
- [161] A. Holzinger, M. Dehmer, F. Emmert-Streib, R. Cucchiara, I. Augenstein, J. Del Ser, W. Samek, I. Jurisica, and N. Díaz-Rodríguez, "Information fusion as an integrative cross-cutting enabler to achieve robust, explainable, and trustworthy medical artificial intelligence," *Information Fusion*, vol. 79, pp. 263–278, 2022.
- [162] K. Yeung, "Recommendation of the council on artificial intelligence (oecd)," *International Legal Materials*, vol. 59, no. 1, pp. 27–34, 2020.
- [163] C. B. Fernandez and P. Hui, "Life, the metaverse and everything: An overview of privacy, ethics, and governance in metaverse," *arXiv preprint arXiv:2204.01480*, 2022.
- [164] L. Floridi, *Ethics, governance, and policies in artificial intelligence*. Springer, 2021.
- [165] J. Qadir, M. Q. Islam, and A. Al-Fuqaha, "Toward accountable human-centered ai: rationale and promising directions," *Journal of Information, Communication and Ethics in Society*, 2022.
- [166] J. Woodward and J. Ruiz, "Analytic review of using augmented reality for situational awareness," *IEEE Transactions on Visualization and Computer Graphics*, 2022.
- [167] U. Ju, L. L. Chuang, and C. Wallraven, "Acoustic cues increase situational awareness in accident situations: A vr car-driving study," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [168] L. Heller and L. Goodman, "What do avatars want now? posthuman embodiment and the technological sublime," in *2016 22nd International Conference on Virtual System & Multimedia (VSMM)*. IEEE, 2016, pp. 1–4.
- [169] B. Shneiderman, "Human-centered artificial intelligence: Reliable, safe & trustworthy," *International Journal of Human-Computer Interaction*, vol. 36, no. 6, pp. 495–504, 2020.
- [170] R. N. Reith, T. Schneider, and O. Tkachenko, "Efficiently stealing your machine learning models," in *Proceedings of the 18th ACM Workshop on Privacy in the Electronic Society*, 2019, pp. 198–210.
- [171] A. Qayyum, K. Ahmad, M. A. Ahsan, A. Al-Fuqaha, and J. Qadir, "Collaborative federated learning for healthcare: Multi-modal COVID-19 diagnosis at the edge," *IEEE Open Journal of the Computer Society*, 2022.
- [172] M. Xu, W. C. Ng, W. Y. B. Lim, J. Kang, Z. Xiong, D. Niyato, Q. Yang, X. S. Shen, and C. Miao, "A full dive into realizing the edge-enabled metaverse: Visions, enabling technologies, and challenges," *arXiv preprint arXiv:2203.05471*, 2022.