# Unleashing the Power of Edge-Cloud Generative AI in Mobile Networks: A Survey of AIGC Services

Minrui Xu, Hongyang Du, Dusit Niyato, *Fellow, IEEE*, Jiawen Kang, Zehui Xiong, Shiwen Mao, *Fellow, IEEE*, Zhu Han, *Fellow, IEEE*, Abbas Jamalipour, *Fellow, IEEE*, Dong In Kim, *Fellow, IEEE*, Xuemin (Sherman) Shen, *Fellow, IEEE*, Victor C. M. Leung, *Life Fellow, IEEE*, and H. Vincent Poor, *Life Fellow, IEEE*

*Abstract*—Artificial Intelligence-Generated Content (AIGC) is an automated method for generating, manipulating, and modifying valuable and diverse data using AI algorithms creatively. This survey paper focuses on the deployment of AIGC applications, e.g., ChatGPT and Dall-E, at mobile edge networks, namely mobile AIGC networks, that provide personalized and customized AIGC services in real time while maintaining user privacy. We begin by introducing the background and fundamentals of generative models and the lifecycle of AIGC services at mobile AIGC networks, which includes data collection, training, fine-tuning, inference, and product management. We then discuss the collaborative cloud-edge-mobile infrastructure and technologies required to support AIGC services and enable users to access AIGC at mobile edge networks. Furthermore, we explore AIGC-driven creative applications and use cases for mobile AIGC networks. Additionally, we discuss the implementation, security, and privacy challenges of deploying mobile AIGC networks. Finally, we highlight some future research directions and open issues for the full realization of mobile AIGC networks.

*Index Terms*—AIGC, Generative AI, Mobile edge networks, Communication and Networking, AI training and inference, Internet technology

M. Xu, H. Du, and D. Niyato are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 308232, Singapore (e-mail: minrui001@e.ntu.edu.sg; hongyang001@e.ntu.edu.sg; dniyato@ntu.edu.sg).

J. Kang is with the School of Automation, Guangdong University of Technology, and Key Laboratory of Intelligent Information Processing and System Integration of IoT, Ministry of Education, Guangzhou 510006, China, and also with Guangdong-HongKong-Macao Joint Laboratory for Smart Discrete Manufacturing, Guangzhou 510006, China (e-mail: kavinkang@gdut.edu.cn).

Z. Xiong is with the Pillar of Information Systems Technology and Design, Singapore University of Technology and Design, Singapore 487372, Singapore (e-mail: zehui_xiong@sutd.edu.sg).

S. Mao is with the Department of Electrical and Computer Engineering, Auburn University, Auburn, AL 36849-5201 USA (email: smao@ieee.org).

Z. Han is with the Department of Electrical and Computer Engineering, University of Houston, Houston, TX 77004 USA, and also with the Department of Computer Science and Engineering, Kyung Hee University, Seoul 446-701, South Korea (e-mail: zhan2@uh.edu).

A. Jamalipour is with the School of Electrical and Information Engineering, University of Sydney, Sydney, NSW 2006, Australia (e-mail: a.jamalipour@ieee.org).

D. I. Kim is with the Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon 16419, South Korea (email:dikim@skku.ac.kr).

X. Shen is with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: sshen@uwaterloo.ca).

V. C. M. Leung is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518061, China, and also with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver BC V6T 1Z4, Canada (E-mail: vleung@ieee.org).

H. V. Poor is with the Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544, USA (e-mail: poor@princeton.edu).
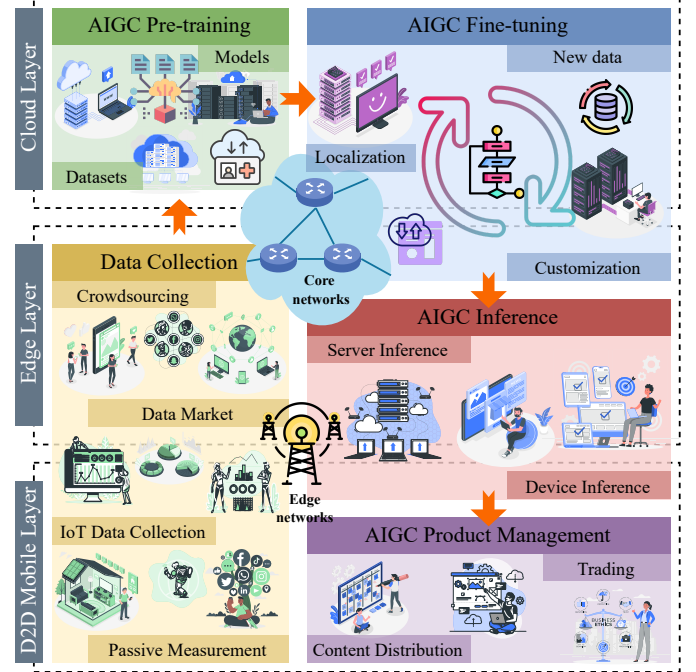
Fig. 1: The overview of mobile AIGC networks, including the cloud layer, the edge layer, and the D2D mobile layer. The lifecycle of AIGC services, including data collection, pre-training, fine-tuning, inference, and product management, is circulated among the core networks and edge networks.

## I. INTRODUCTION

### A. Background

In recent years, artificial intelligence-generated content (AIGC) has emerged as a novel approach to the production, manipulation, and modification of data. By utilizing AI technologies, AIGC automates content generation alongside traditionally professionally-generated content (PGC) and user-generated content (UGC) [1]–[3]. With the marginal cost of data creation reduced to nearly zero, AIGC, e.g., ChatGPT, promises to supply a vast amount of synthetic data for AI development and the digital economy, offering significant productivity and economic value to society. The rapid growth of AIGC capabilities is driven by the continuous advancements in AI technology, particularly in the areas of large-scale and multimodal models [4], [5]. A prime example of this progress is the development of DALL-E [6], an AI system based

on OpenAI's state-of-the-art GPT-3 language model, which consists of 175 billion parameters and is designed to generate images by predicting successive pixels. In its latest iteration, DALL-E2 [7], a diffusion model is employed to reduce noise generated during the training process, leading to more refined and novel image generation. In the context of text-to-image generation using AIGC models, the language model serves as a guide, enhancing semantic coherence between the input prompt and the resulting image. Simultaneously, the AIGC model processes existing image attributes and components, generating limitless synthesis images from existing datasets.

Based on large-scale pre-trained models with billions of parameters, AIGC services are designed to enhance knowledge and creative work fields that employ billions of people. By leveraging generative AI, these fields can achieve at least a 10% increase in efficiency for content creation, potentially generating trillions of dollars in economic value [8]. AIGC can be applied to various forms of text generation, ranging from practical applications, such as customer service inquiries and messages, to creative tasks like activity tracking and marketing copywriting [9]. For example, OpenAI's ChatGPT [10] can automate the generation of socially valuable content based on user-provided prompts. Through extended and coherent conversations with ChatGPT, individuals from diverse professions from all walks of life, can seek assistance in debugging code, discovering healthy recipes, writing scripts, and devising marketing campaigns. In the realm of image generation, AIGC models can process existing images according to their attributes and components, enabling end-to-end image synthesis, such as generating complete images directly from existing ones [7]. Moreover, AIGC models hold immense potential for cross-modal generation, as they can spatially process existing video attributes and simultaneously process multiple video clips automatically [11].

The benefits of AIGC in content creation, when compared to PGC and UGC, are already apparent to the public. Specifically, generative AI models can produce high-quality content within seconds and deliver personalized content tailored to users' needs [2]. Over time, the performance of AIGC has significantly improved, driven by enhanced models, increased data availability, and greater computational power [12]. On one hand, superior models [4], such as diffusion models, have been developed to provide more robust tools for cross-modal AIGC generation. These advancements are attributed to the foundational research in generative AI models and the continuous refinement of learning paradigms and network structures within generative deep neural networks (DNNs). On the other hand, data and computing power for generative AI training and inference have become more accessible as networks grow increasingly interconnected [9], [13]. For instance, AIGC models that require thousands of GPUs can be trained and executed in cloud data centers, enabling users to submit frequent data generation requests over core networks.

### B. Motivation

Although AIGC is acknowledged for its potential to revolutionize existing production processes, users accessing AIGC services on mobile devices currently lack support for interactive and resource-intensive data generation services [14], [25]. Initially, the robust computing capabilities of cloud data centers can be utilized to train AIGC pre-training models, such as GPT-3 for ChatGPT and GPT-4 for ChatGPT Plus. Subsequently, users can access cloud-based AIGC services via the core network by executing AIGC models on cloud servers. However, due to their remote nature, cloud services exhibit high latency. Consequently, deploying interaction-intensive AIGC services on mobile edge networks, i.e., mobile AIGC networks, as shown in Fig. 1, should be considered a more practical option [26]–[28]. In detail, the motivations for developing mobile AIGC networks include

- *Low-latency:* Instead of directing requests for AIGC services to cloud servers within the core network, users can access low-latency services in mobile AIGC networks [29]. For example, users can obtain AIGC services directly in radio access networks (RANs) by downloading pre-trained models to edge servers and mobile devices for fine-tuning and inference, thereby supporting real-time, interactive AIGC.
- *Localization and Mobility:* In mobile AIGC networks, base stations with computing servers at the network's edge can fine-tune pre-trained models by localizing service requests [30], [31]. Furthermore, users' locations can serve as input for AIGC fine-tuning and inference, addressing specific geographical demands. Additionally, user mobility can be integrated into the AIGC service provisioning process, enabling dynamic and reliable AIGC service provisioning.
- *Customization and Personalization:* Local edge servers can adapt to local user requirements and allow users to request personalized services based on their preferences while providing customized services according to local service environments. On one hand, edge servers can tailor AIGC services to the needs of the local user community by fine-tuning them accordingly [2]. On the other hand, users can request personalized services from edge servers by specifying their preferences.
- *Privacy and Security:* AIGC users only need to submit service requests to edge servers, rather than sending preferences to cloud servers within the core network. Therefore, the privacy and security of AIGC users can be preserved during the provisioning, including fine-tuning and inference, of AIGC services.

As illustrated in Fig. 1, when users access AIGC services on mobile edge networks through edge servers and mobile devices, limited computing, communication, and storage resources pose challenges for delivering interactive and resource-intensive AIGC services. First, resource allocation on edge servers must balance the tradeoff among accuracy, latency, and energy consumption of AIGC services at edge servers. In addition, computationally intensive AIGC tasks can be offloaded from mobile devices to edge servers, improving inference latency and service reliability. Moreover, AI models that generate content can be cached in edge networks, similar to content delivery networks (CDNs) [32], [33], to

TABLE I: Summary of related works versus our survey.

| Year | Ref. | Contributions | AIGC Algorithms | AIGC Applications | Edge Intelligence |
|------|------|---------------|-----------------|-------------------|-------------------|
| 2019 | [14] | Introduce mobile edge intelligence, and discuss the infrastructure, implementation methodologies, and use cases | ✗ | ✗ | ✓ |
| 2020 | [15] | Present the implementation challenges of federated learning at mobile edge networks | ✓ | ✗ | ✓ |
| | [12] | Discuss the visions, implementation details, and applications of the convergence of edge computing and DL | ✓ | ✗ | ✓ |
| 2021 | [16] | Investigate the copyright laws regarding AI-generated music | ✓ | ✓ | ✗ |
| | [1] | Illustrate the interaction of art and AI from two perspectives, i.e., AI for art analysis and AI for art creation | ✗ | ✓ | ✗ |
| | [2] | Discuss the application of computational arts in Metaverse to create surrealistic cyberspace | ✓ | ✓ | ✗ |
| | [17] | Investigate the deployment of distributed learning in wireless networks | ✗ | ✗ | ✓ |
| | [18] | Provide a comprehensive overview of the major approaches, datasets, and metrics used to synthesize and process multimodal images | ✓ | ✓ | ✗ |
| | [19] | Propose a novel conceptual architecture for 6G networks, which consists of holistic network virtualization and pervasive network intelligence | ✗ | ✗ | ✓ |
| 2022 | [20] | Discusses the visions and potentials of low-power, low-latency, reliable, and trustworthy edge intelligence for 6G wireless networks | ✗ | ✗ | ✓ |
| | [4] | Provide comprehensive guidance and comparison among advanced generative models, including GAN, energy-based models, VAE, autoregressive models, flow-based models, and diffusion models | ✓ | ✗ | ✗ |
| | [21] | Present fundamental algorithms, classification and applications of diffusion models | ✓ | ✗ | ✗ |
| | [9] | Provide a comprehensive overview of generation and detection methods for machine-generated text | ✓ | ✓ | ✗ |
| | [22] | Provide a comprehensive examination of what, why, and how edge intelligence and blockchain can be integrated | ✗ | ✗ | ✓ |
| | [23] | Introduce the architecture of edge-enabled Metaverse and discuss enabling technologies in communication, computing, and blockchain | ✗ | ✓ | ✓ |
| 2023 | [24] | Summarize existing works on the generation of gestures with simultaneous speeches based on deep generative models | ✓ | ✓ | ✗ |
| | Ours | Investigate the deployment of mobile AIGC networks via collaborative cloud-edge-mobile infrastructure, discuss creative mobile applications and exemplary use cases, and identify existing implementation challenges | ✓ | ✓ | ✓ |

minimize delays in accessing the model. Finally, mobility management and incentive mechanisms should be explored to encourage user participation in both space and time. Compared to traditional AI, AIGC technology requires overall technical maturity, transparency, robustness, impartiality, and insightfulness of the algorithm for effective application implementation. From a sustainability perspective, AIGC can use both existing and synthetic datasets as raw materials for generating new data. However, when biased data are used as raw data, these biases persist in the knowledge of the model, which inevitably leads to unfair results of the algorithm. Finally, static AIGC models rely primarily on templates to generate machine-generated content that may have similar text and output structures.

## C. Related Works and Contributions

In this survey, we provide an overview of research activities related to AIGC and mobile edge intelligence, as illustrated in Fig. 2. Given the increasing interest in AIGC, several surveys on related topics have recently been published. Table I presents a comparison of these surveys with this paper.

The study in [34] provides a comprehensive overview of the current AIGC models published by researchers and the industry. The authors identify nine categories summarizing the evolution of generative AI models, including text-to-text, text-to-image, text-to-audio, text-to-video, text-to-3D, text-to-code, text-to-science, image-to-text, and other models. In addition, they reveal that only six organizations with enormous computing power and highly skilled and experienced teams can deploy these state-of-the-art models, which is even fewer than the number of categories. Following the taxonomy of generative AI models developed in [34], other surveys discuss generative AI models in detail subsequently. The study in [9] examines existing methods for generating text and detecting models. The study in [18] provides a comprehensive overview of the major approaches, datasets, and evaluation metrics
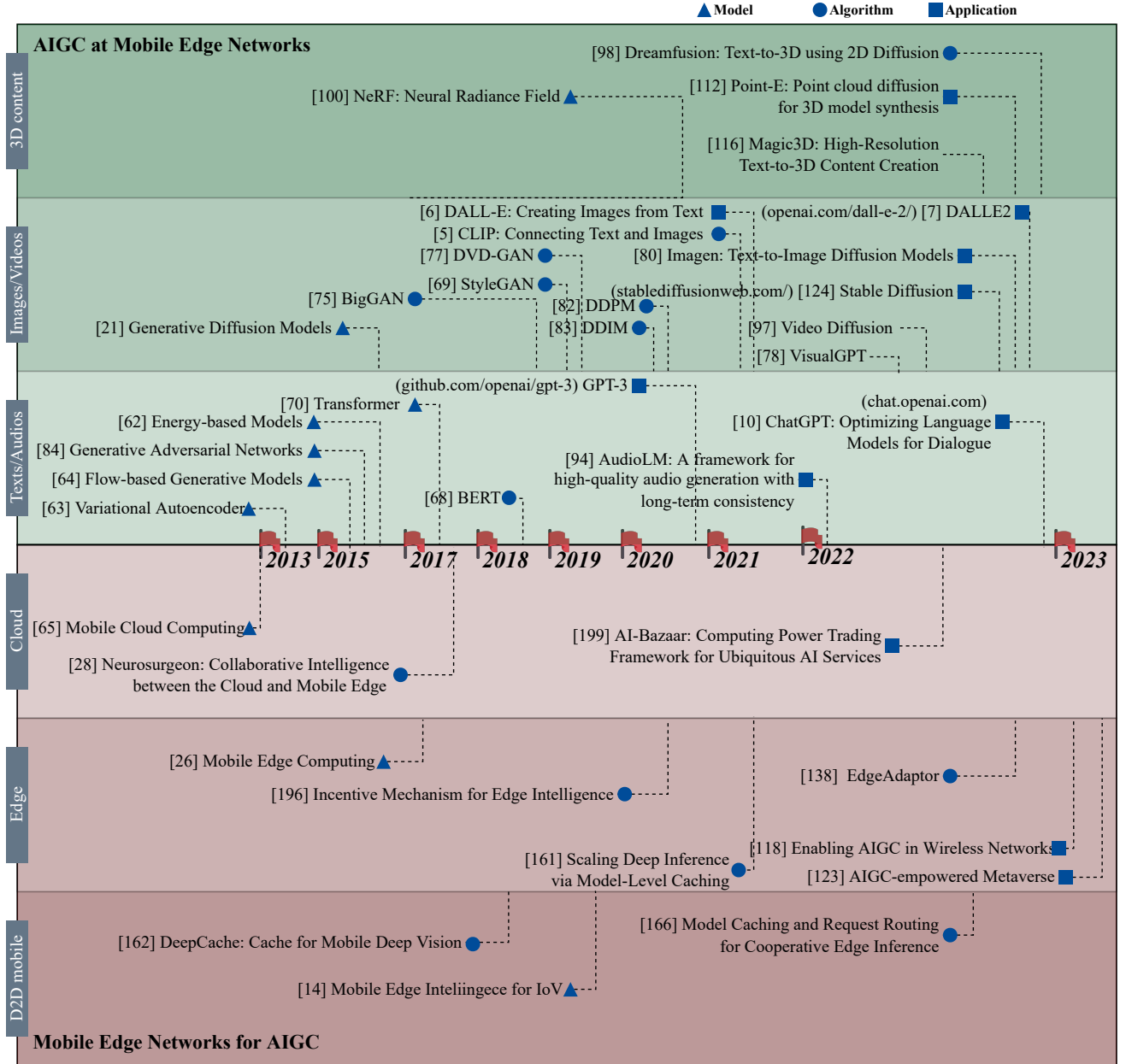
Fig. 2: The development roadmap of AIGC and mobile edge networks from 2013 to Jan 2023. From the perspective of AIGC technology development, AIGC has evolved from generating text and audio to generating 3D content. From the perspective of mobile edge computing, computing has gradually shifted from cloud data centers to D2D mobile computing.

for multimodal image synthesis and processing. Based on techniques of speech and image synthesis, the study in [24] summarizes existing works on the generation of gestures with simultaneous speeches based on deep generative models. The study in [16] investigates the copyright laws regarding AI-generated music, which includes the complicated interactions among AI tools, developers, users, and the public domain. The study in [4] provides comprehensive guidance and comparison among advanced generative models, including GANs, energy-based models, variational autoencoder (VAE), autoregressive models, flow-based models, and diffusion models. As diffusion models draw tremendous attention in generating creative data, the study in [21] presents fundamental algorithms and

comprehensive classification for diffusion models. Based on these algorithms, the authors [1] illustrate the interaction of art and AI from two perspectives, i.e., AI for art analysis and AI for art creation. In addition, the authors in [2] discuss the application of computational arts in the Metaverse to create surrealistic cyberspace.

In 6G [19], mobile edge intelligence based on edge computing systems, including edge caching, edge computing, and edge intelligence, for intelligent mobile networks, is introduced in [14]. The study in [17] investigates the deployment of distributed learning in wireless networks. The study [15] provides a guide to federated learning and a comprehensive overview of implementing Federated Learning (FL) at mobile

Fig. 3: The outline of this survey, where we introduce the provisioning of AIGC services at mobile edge networks and highlight some essential implementation challenges about mobile edge networks for provisioning AIGC services.

edge networks. The authors offer a detailed analysis of the challenges of implementing FL, including communication costs, resource allocation, privacy, and security. In [12], various application scenarios and technologies for edge intelligence and intelligent edges are presented and discussed in detail. In addition, the study [20] discusses the visions and potentials of low-power, low-latency, reliable, and trustworthy edge intelligence for 6G wireless networks. The study [22] explores how blockchain technologies can be used to enable edge intelligence and how edge intelligence can support the deployment of blockchain at mobile edge networks. The authors provide a comprehensive review of blockchain-driven edge intelligence, edge intelligence-amicable blockchain, and their implementation at mobile edge networks. We also [23] provide a vision of realizing the Metaverse at mobile edge

networks. In detail, enabling technologies and challenges are discussed, including communication and networking, computing, and blockchain.

Distinct from existing surveys and tutorials, our survey concentrates on the deployment of mobile AIGC networks for real-time and privacy-preserving AIGC service provisioning. We introduce the current development of AIGC and collaborative infrastructure in mobile edge networks. Subsequently, we present the technologies of deep generative models and the workflow of provisioning AIGC services within mobile AIGC networks. Additionally, we showcase creative applications and several exemplary use cases. Furthermore, we identify implementation challenges, ranging from resource allocation to security and privacy, for the deployment of mobile AIGC networks. The *contributions of our survey* are as follows.

- We initially offer a tutorial that establishes the definition, lifecycle, models, and metrics of AIGC services. Then, we propose the mobile AIGC networks, i.e., provisioning AIGC services at mobile edge networks with collaborative mobile-edge-cloud communication, computing, and storage infrastructure.
- We present several use cases in mobile AIGC networks, encompassing creative AIGC applications for text, images, video, and 3D content generation. We summarize the advantages of constructing mobile AIGC networks based on these use cases.
- We identify crucial implementation challenges in the path to realizing mobile AIGC networks. The implementation challenges of mobile AIGC networks stem not only from dynamic channel conditions but also from the presence of meaningless content, insecure content precepts, and privacy leaks in AIGC services.
- Lastly, we discuss future research directions and open issues from the perspectives of networking and computing, machine learning (ML), and practical implementation considerations, respectively.

As the outline illustrated in Fig. 3, the survey is organized as follows. Section II examines the background and fundamentals of AIGC. Section III presents the technologies and collaborative infrastructure of mobile AIGC networks. The applications and advantages of mobile AIGC networks are discussed in Section IV, and potential use cases are shown in Section V. Section VI addresses the implementation challenges. Section VII explores future research directions. Section VIII provides the conclusions.

## II. BACKGROUND AND FUNDAMENTALS OF AIGC

In this section, the background and fundamentals of AIGC technology are presented in this section. Specifically, we examine the definition of AIGC, its classification, and the technological lifecycle of AIGC in mobile networks. Finally, we introduce ChatGPT as a use case, which is the most famous and revolutionary application of AIGC.

### A. Definitions of PGC, UGC, and AIGC

In the next generation of the Internet, i.e. Web 3.0 and Metaverse [35], there are three primary forms of content [1], including PGC, UGC, and AIGC.

*1) Professionally-generated Content:* PGC refers to professional-generated digital content [36]. Here, the generators are individuals or organizations with professional skills, knowledge, and experience in a particular field, e.g., journalists, editors, and designers. As these experts who create PGC are typically efficient and use specialized tools, PGC has the advantages in terms of *automation* and *multimodality*. However, because PGC is purposeful, the *diversity* and *creativity* of PGC can be limited.

*2) User-generated Content:* UGC refers to digital material generated by users, rather than by experts or organizations [37]. The users include website visitors and social media users. UGC can be presented in any format, including text, photos, video, and audio. The barrier for users to creating UGC

is being lowered. For example, some websites[1] allow users to create images with a high degree of freedom on a pixel-by-pixel basis. As a result, UGC is more *creative* and *diverse*, thanks to a wide user base. However, UGC is less *automated* and less *multimodal* than the PGC that is generated by experts.

*3) AIGC:* AIGC is generated by using generative AI models according to input from users. Because AI models can learn the features and patterns of input data from the human artistic mind, they can develop a wide range of content. The recent success of text-to-image applications based on the diffusion model [38] and the ChatGPT based on transformer [10] has led to AIGC gaining a lot of attention. We have defined the AIGC according to its characteristics as follows

- Automatic: AIGC is generated by AI models automatically. After the AI model has been trained, users only need to provide input, such as the task description, to efficiently obtain the generated content. The process, from input to output, does not require user involvement and is done automatically by the AI models.
- Creativity: AIGC refers to an idea or item that is innovative. For example, AIGC is believed to be leading to the development of a new profession, called Prompt Engineer [39], which aims to improve human interaction with AI. In this context, the prompt serves as the starting point for the AI model, and it significantly impacts the originality and quality of the generated content. A well-crafted prompt that is precise and specific results in more relevant and creative content than a vague or general prompt.
- Multimodal: The AI models to generate AIGC can handle multimodal input and output. For example, ChatGPT [10] allows conversational services that employ text as input and output, DALL-E 2 [40] can create original, realistic images from a text description, and AIGC services with voice and 3D models as input or output are progressing [41].
- Diverse: AIGC is diverse in service personalization and customization. On the one hand, users can adjust the input to the AI model to suit their preferences and needs, resulting in a personalized output. On the other hand, AI models are trained to provide diverse outputs. For example, consider the DALL-E 2 as an example, the model can generate images of individuals that more correctly represent the diversity of the global population, even with the same text input.
- Extendedly valuable: AIGC should be extendedly valuable to society, economics, and humanity [42]. For example, AI models can be trained to write medical reports and interpret medical images, enabling healthcare personnel to make accurate diagnoses.

AIGC provides various advantages over PGC and UGC, including better efficiency, originality, diversity, and flexibility. The reason is that AI models can produce vast amounts of material quickly and develop original content based on established patterns and principles. These advantages have

---

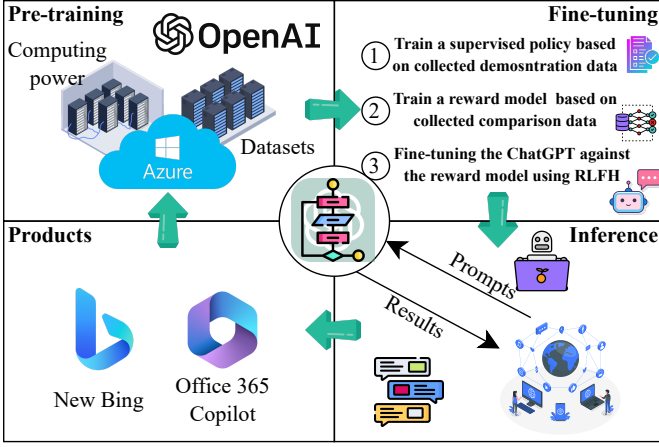[1]Example of a website that allows users to create their own UGC: https://ugc-nft.io/Home

Fig. 4: The four development stages of ChatGPT, including pre-training, fine-tuning, inference, and product management.

led to the growing creative applications of the AIGC models, which are discussed in Section IV-A1.

### B. Serving ChatGPT at Mobile Edge Networks

ChatGPT, developed by OpenAI, excels at generating human-like text and engaging in conversations [10]. Based on the GPT-3 [43], this transformer-based neural network model can produce remarkably coherent and contextually appropriate text. Among its primary advantages, ChatGPT is capable of answering questions, providing explanations, and assisting with various tasks in a manner nearly indistinguishable from human responses. As illustrated in Fig. 4, the development of ChatGPT involves four main stages, including pre-training, fine-tuning, inference, and product management.

*1) Pre-training:* In the initial stage, known as pre-training, the foundation model of ChatGPT, GPT-3, is trained on a large corpus of text, which includes books, articles, and other information sources. This process enables the model to acquire knowledge of language patterns and structures, as well as the relationships between words and phrases. The base model, GPT-3, is an autoregressive language model with a Transformer architecture that has 175 billion parameters, making it one of the largest language models available. During pre-training, GPT-3 is fed with a large corpus of text from diverse sources, such as books, articles, and websites for self-supervised learning, where the model learns to predict the next word in a sentence given the context. To train the foundation model, the technique used is called maximum likelihood estimation, where the model aims to maximize the probability of predicting the next word correctly. Training GPT-3 demands significant computational resources and time, typically involving specialized hardware like graphics processing units (GPUs) or tensor processing units (TPUs). The exact resources and time required to depend on factors such as model size, dataset size, and optimization techniques.

*2) Fine-tuning:* The fine-tuning stage of ChatGPT involves adapting the model to a specific task or domain, such as customer service or technical support, in order to enhance its accuracy and relevance for that task. To transform ChatGPT

into a conversational AI, a supervised learning process is employed using a dataset containing dialogues between humans and AI models [44]. To optimize ChatGPT's parameters, a reward model for reinforcement learning is built by ranking multiple model responses by quality. Alternative completions are ranked by AI trainers, and the model uses these rankings to improve its performance through several iterations of Proximal Policy Optimization [45]. This technique allows ChatGPT to learn from its mistakes and improve its responses over time.

*3) Inference:* In the inference stage, ChatGPT generates text based on a given input or prompt, testing the model's ability to produce coherent and contextually appropriate responses relevant to the input. ChatGPT generates responses by leveraging the knowledge it acquired during pre-training and fine-tuning, analyzing the context of the input to generate relevant and coherent responses. In-context learning involves analyzing the entire context of the input [46], including the dialogue history and user profile, to generate responses that are personalized and tailored to the user's needs. ChatGPT employs chain-of-thought to generate responses that are coherent and logical, ensuring that the generated text is not only contextually appropriate but also follows a logical flow. The resources consumed during inference are typically much lower than those required for training, making real-time applications and services based on ChatGPT computationally feasible.

*4) Product Management:* The final product management phase involves deploying the model in a production environment and ensuring its smooth and efficient operation. In the context of mobile edge networks, the applications of AI-powered tools such as the new Bing [47] and Office 365 Copilot [48] could be particularly useful due to their ability to provide personalized and contextually appropriate responses while conserving resources. The new Bing offers a new type of search experience with AI-powered features such as detailed replies to complex questions, summarized answers, and personalized responses to follow-up questions, while Office 365 Copilot, powered by GPT-4 from OpenAI, provides assistance with generating documents, emails, presentations, and other tasks in Microsoft 365 apps and services. These tools can be integrated into mobile edge networks with specialized techniques that balance performance and accuracy while preserving data integrity.

- New bing: The new Bing offers a set of AI-powered features that provide a new type of search experience, including detailed replies to complex questions, summarized answers, and personalized responses to follow-up questions. Bing also offers creative tools such as assistance with writing poems and stories. In the context of mobile edge networks, Bing's ability to consolidate reliable sources across the web and provide a single, summarized answer could be particularly useful for users with limited resources. Additionally, Bing's ability to generate personalized responses based on user behavior and preferences could improve the experience of users in mobile edge networks.
- Office 365 copilot: Microsoft has recently launched an AI-powered assistant named Office 365 Copilot, which can be summoned from the sidebar of Microsoft 365 apps

and services. Copilot can help users generate documents, emails, and presentations, as well as provide assistance with features such as PivotTables in Excel. It can also transcribe meetings, remind users of missed items, and provide summaries of action items. However, when deploying Copilot in mobile edge networks, it is important to keep in mind the limited resources of these devices and to develop specialized techniques that can balance performance and accuracy while preserving data integrity.

In addition to the previously mentioned commercial applications, ChatGPT holds substantial commercial potential owing to its capacity for producing human-like text, which is characteristically coherent, pertinent, and contextually fitting. This language model can be fine-tuned to accommodate a diverse array of tasks and domains, rendering it highly adaptable for numerous applications. ChatGPT exhibits remarkable proficiency in comprehending and generating text across multiple languages. Consequently, it can facilitate various undertakings, such as composing emails, developing code, generating the content, and offering explanations, ultimately leading to enhanced productivity. By automating an assortment of tasks and augmenting human capabilities, ChatGPT contributes to a paradigm shift in the nature of human work, fostering new opportunities and revolutionizing industries. In addition to ChatGPT, more use cases developed by various generative AI models are discussed in Section V.

### C. Life-cycle of AIGC at Mobile Edge Networks

AIGC has gained tremendous attention as a technology superior to PGC and UGC. However, the lifecycle of the AIGC is also more elaborate. In the following, we discuss the AIGC lifecycle with mobile edge network enablement:

*1) Data Collection:* Data collection is an integral component of AIGC and plays a significant role in defining the quality and diversity of the material created by AI systems [49]. The data used to train AI models influences the patterns and relationships that the AI models learn and, consequently, the output. There are several data collection techniques for AIGC:

- Crowdsourcing: Crowdsourcing is the process of acquiring information from a large number of individuals, generally via the use of online platforms [50]. Crowdsourced data may be used to train ML models for text and image generation, among other applications. One common example is the use of Amazon Mechanical Turk[2], where individuals are paid to perform tasks such as annotating text or images, which can then be used to train AIGC models.
- Data Market: Another way to obtain data is to buy it from a data provider. For example, Datatang[3] is a firm that offers high-quality datasets and customized data services to assist businesses in enhancing the performance of their AI models. By giving access to varied, high-quality data, Datatang enables organizations to train AI models that

are more accurate and effective, resulting in enhanced business performance and results.

- Internet-of-Things (IoT) data collection: In IoT, edge devices can help to collect the data, e.g., Global Positioning System (GPS) records and wireless sensing data [51]. For example, installing sensors in mobile phones that can track the movement and location of the devices or users [52]. The sensors can be used to collect data on the location, speed, and direction of movement of the device. These data are important for the implementation of personalized AIGC models.
- Passive measurement: Passive data collection can be achieved with the help of edge networks [53]. In the smart city, sensors can be placed at strategic locations, such as on lamp posts, buildings, or other structures, to collect data on various aspects of the city environment. The data obtained by the sensors might be used to train AI models, which could subsequently be utilized to produce insights on air quality, traffic flow, and pedestrian density. Using data obtained from air quality sensors, for instance, an AI model may be trained to forecast air quality. The model may then be used to create a real-time map of the city's air quality, which could be used to guide policy choices about the management of air quality.

After the data has been collected, the data is then used to train the AIGC model.

*2) Pre-training:* The collected data is used to train the AIGC model. In mobile networks, training is typically done by central servers with powerful computing power. During the training process, the generative model automatically learns the patterns and features in the data and predicts the target outcome. We introduce several generative AI technologies in Section III-B, including Generative Adversarial Networks (GANs), VAE, Flow-based models, and diffusion models. These different training techniques have different strengths and weaknesses. The choice of technique depends on the specific requirements of the AIGC task, the available data, the desired output, and the computational resources available. After training is complete, cloud data centers can accept requests uploaded by network users to perform subsequent fine-tuning and inference tasks. Alternatively, cloud data centers can deliver the trained AIGC models down to network edge servers, which can process user requests locally.

*3) Fine-tuning:* Fine-tuning in AIGC is the process of adjusting a pre-trained AIGC model to new tasks or domains by including a modest quantity of extra data. This approach can be used to enhance the model's performance on a given task or in a specific area by adjusting the AI model's parameters to suit the new data better. In mobile networks, tasks of fine-tuning can be performed by the edge network, using the small-size dataset uploaded by mobile users.

*4) Inference:* Using the trained AIGC model, inference can be done, which involves generating the desired content based on the input. AIGC models are traditionally managed via centralized servers, such as the Hugging Face platform [54]. In this setting, a large number of users make requests to the central server, wait in line, and obtain the requested services. Researchers aim to install AIGC services on edge networks

---

[2]The website of Amazon Mechanical Turk as a crowdsourcing marketplace: https://www.mturk.com/

[3]The website of Datatang: https://www.datatang.ai/

to prevent request congestion and optimize service latency. Edge devices have the sufficient computational capacity for AIGC inference and are closer to consumers than central servers. Therefore, users can interact with devices with a reduced transmission delay. In addition, as AIGC services are dispersed to several edge devices, the latency can be significantly reduced.

*5) Product Management:* The preceding stages cover the AIGC generation. However, as an irreplaceable online property comparable to NFT, AIGC possesses unique ownership, copyright, and worth for each content. Consequently, the preservation and management of AIGC products should be incorporated into the AIGC life cycle. Specifically, we refer to the party requesting the production of the AIGC as producers, e.g., mobile users or companies, who hire AIGC generators, e.g., network servers, to perform the AIGC tasks. Then, the main process in AIGC product management includes:

- *Distribution:* After the content is generated in network edge servers, the producers acquire ownership of the AIGC products. Consequently, they have the right to distribute these products to social media or AIGC platforms through edge networks
- *Trading:* Since AIGC products are regarded as a novel kind of non-fungible digital properties, they can be traded. The trading process can be modeled as a fund ownership exchange between two parties.

To implement the aforementioned AIGC lifecycle in mobile networks, we further investigate the technical implementation of AIGC in the following section.

## III. TECHNOLOGIES AND COLLABORATIVE INFRASTRUCTURE OF MOBILE AIGC NETWORKS

In this section, we delve into the technologies and collaborative infrastructure of mobile AIGC networks. This section aims to provide a comprehensive understanding of the rationale and objectives of edge computing systems designed to support AIGC. Before we explore the design of these systems, it is crucial to establish the performance metrics that measure whether the system can maximize user satisfaction and utility.

### A. Evaluation Metrics of Generative AI Models and Services

We first discuss several metrics for assessing the quality of AIGC models, which can be used by AIGC service providers and users in mobile networks.

*1) Inception Score:* The Inception Score (IS) can be used to measure the accuracy of images generated by AIGC models in the mobile network [55]. The IS is based on the premise that high-fidelity generated images should have high-class probabilities, which suggest a reliable classification model, and a low Kullback-Leibler (KL) divergence between the projected class probability and a reference class distribution. To compute the IS, an exponential function is applied to the KL divergence between the anticipated class probabilities and the reference class distribution. The resulting value is then averaged over all created photos to obtain the IS. A higher IS indicates better overall image quality.

*2) Frechet Inception Distance:* The Frechet Inception Distance (FID) has emerged as a well-established metric for evaluating the effectiveness of generative models, particularly GANs, in terms of image quality and diversity [56]. FID leverages a pre-trained Inception network to calculate the distance between actual and synthetic image embeddings. This metric can be used by AIGC model providers to evaluate the quality of their generative models in mobile networks. Additionally, users can assess the capabilities of AIGC service providers through multiple requests for services based on FID measurements. However, when evaluating conditional text-to-image synthesis, FID only measures the visual quality of the output images, ignoring the adequacy of their conditioning on the input text [57]. Thus, while FID is an excellent evaluation metric for assessing image quality and diversity, it is limited when applied to conditional text-to-image synthesis.

*3) R-Precision:* R-Precision is a standard metric to evaluate how AI-generated images align with text inputs [58]. In mobile networks, the AIGC model producers can retrieve matching text from 100 text candidates using the AI-generated image as a query. The R-Precision measures the proportion of relevant items retrieved among the top-R retrieved items, where R is typically set to 1. Specifically, the Deep Attentional Multimodal Similarity Model (DAMSM) is commonly used to compute the text-image retrieval similarity score [59]. DAMSM maps each subregion of an image and its corresponding word in the sentence to a joint embedding space, allowing for the measurement of fine-grained image-text similarity for retrieval. However, it should be noted that text-to-image AIGC models can directly optimize the DAMSM module used to calculate R-Precision. This results in the metric being model-specific and less objective, limiting the evaluation of AIGC models in mobile networks.

*4) CLIP-R-Precision:* CLIP-R-Precision is an assessment metric to address the model-specific character of the R-Precision metric [60]. Instead of the conventional DAMSM, the suggested measure uses the latest multimodal CLIP model [5] to obtain R-Precision scores. Here, CLIP is trained on a massive corpus of web-based image-caption pairings and is capable, via a contrastive aim, of bringing together the two embeddings (visual and linguistic). Thus, the CLIP-R-Precision can provide a more objective evaluation of text-to-image AIGC model performance in mobile networks.

*5) Quality of Experience:* The Quality of Experience (QoE) metric plays a critical role in evaluating the performance of AIGC in mobile network applications. QoE measures user satisfaction with the generated content, considering factors such as visual quality, relevancy, and utility. Gathering and analyzing user surveys, interaction, and behavioral data are standard methods used to determine QoE. In addition, the definition of QoE can vary depending on the objectives of the mobile network system designer and the user group being considered. With the aid of QoE, AIGC performance can be improved, and new models can be created to meet user expectations. It is essential to account for QoE when analyzing the performance of AIGC in mobile network applications to ensure that the generated content meets user expectations and provides a great user experience.
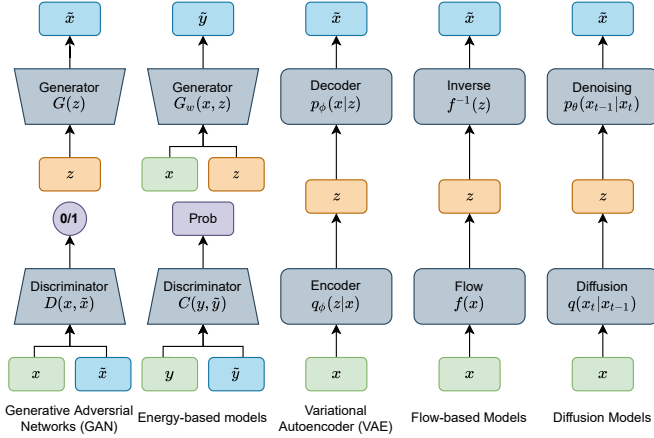
Fig. 5: The model architecture of generative AI models, including generative adversarial networks, energy-based models, variational autoencoder, flow-based models, and diffusion models.

Based on the aforementioned evaluation metrics, diverse and valuable synthetic data can be generated from deep generative models. Therefore, in the next section, we introduce several generative AI models for mobile AIGC networks.

### B. Generative AI Models

The objective of a generative AI model is to fit the true data distribution of input data by iterative training. Users can generate novel data using this approximate model as the model can fit into the distribution. As shown in Fig. 5, this section mainly introduces five basic generative models, including GANs, energy-based models, VAE, flow-based models, and diffusion models.

*1) Generative Adversarial Networks:* The GAN [61] is a fundamental framework for AIGC, comprising a generative model and a discriminative model. The generative network aims to generate data that is as realistic and similar to the original data as possible to deceive the discriminative model, based on the data in the original dataset. Conversely, the discriminant model's task is to differentiate between real and fake instances. During the GAN training process, the two networks continually enhance their performance by competing against each other until they reach a stable equilibrium. By the end of the training process, the discriminator network is no longer able to differentiate between real and fake data. However, GANs have limited control over the output and can produce meaningless images. Moreover, they generate low-resolution images, only augment the existing dataset rather than creating new content on the original dataset, and cannot generate new content across modalities.

*2) Energy-based Generative Models:* Energy-based generative models [62] are probabilistic generative models that represent input data using energy values and model the data by minimizing these values. The energy-based models function by defining an energy function and then minimizing the energy value of the input data through optimization and training. This

approach has the advantage of being easily comprehensible, and the models exhibit excellent flexibility and generalization ability in providing AIGC services.

*3) Variational Autoencoder:* The VAE [63] consists of two main components: an encoder and a decoder network. The encoder converts the input data into the mean and variance of the latent measures and uses these parameters to sample the latent space and generate the latent measures. The decoder takes the latent variables as input and generates new data. The data reconstruction and data generation tasks can be accomplished by training the encoder and decoder together. Unlike GANs, which are trained using a supervised learning approach, VAE uses an unsupervised learning approach. Thus, the VAE generates data by sampling from the learned distribution, while the GAN generates data by approximating the data distribution using the generator network.

*4) Flow-based Generative Models:* Flow-based generative models [64] facilitate the data generation process by employing probabilistic flow formulations. Additionally, these models compute gradients during generation using backpropagation algorithms, enhancing training and learning efficiency. Consequently, flow-based models in mobile edge networks present several benefits. One such advantage is computational efficiency. Flow-based models can directly compute the probability density function during generation, circumventing resource-intensive calculations. This promotes more efficient computation within mobile edge networks.

*5) Generative Diffusion Models:* Diffusion models are likelihood-based models trained with Maximum Likelihood Estimation (MLE) [21], as opposed to GANs trained with a minimax game between the generator and the discriminator. Therefore, the pattern collapses and thus the training instabilities can be avoided. Specifically, diffusion models are inspired by non-equilibrium thermodynamics theory. They learn the inverse diffusion process to construct the desired data sample from noise by defining a Markov chain of diffusion steps that gradually add random noise to the data. In addition, diffusion can mathematically transform the computational space of the model from pixel space to a low-dimensional space called latent space. This reduces the computational cost and time required and improves the training efficiency of the model. Unlike VAE or flow-based models, diffusion models are learned using a fixed procedure, and the hidden variables have high dimensions that are the same as the original data.

### C. Collaborative Infrastructure for Mobile AIGC Networks

By asking ChatGPT the question "Integrating AI-generated content and mobile edge networks, please define mobile AIGC networks in one sentence," we can get the answer "*Mobile AIGC networks are a fusion of AI-generated content and mobile edge networks, enabling rapid content creation, delivery, and processing at the network's edge for enhanced user experiences and reduced latency.*" (from Mar. 14 Version based on GPT-4) To support the pre-training, fine-tuning, and inference of the aforementioned models, substantial computation, communication, and storage resources are necessary. Consequently, to provide low-latency and personalized AIGC
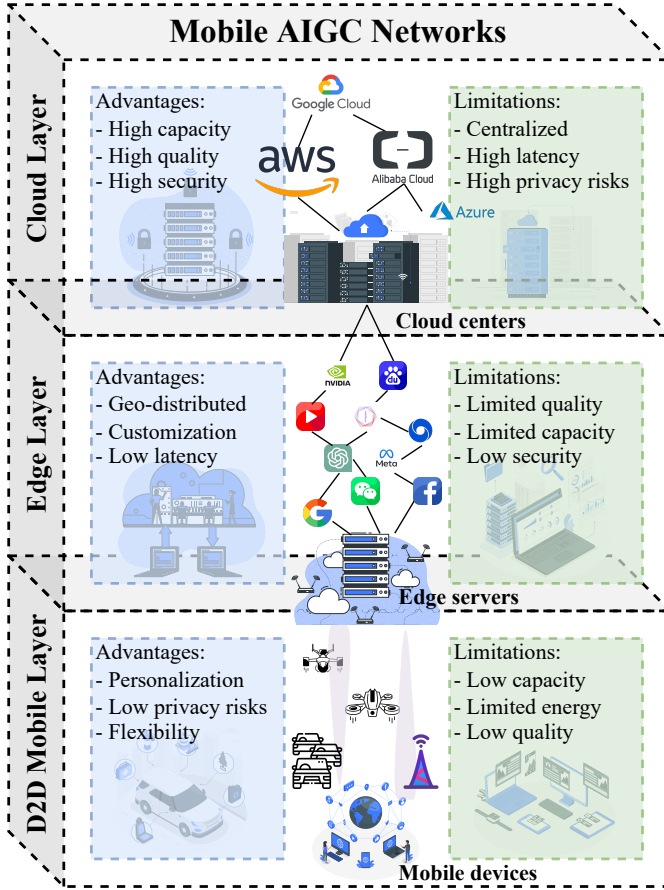
Fig. 6: The collaborative cloud-edge-mobile infrastructure for mobile AIGC networks. The advantages and limitations of provisioning AIGC services in each layer are elaborated.

services, a collaborative cloud-edge-mobile AIGC framework shown in Fig. 6 is essential, requiring extensive cooperation among heterogeneous resource shareholders.

*1) Cloud Computing:* In mobile AIGC networks, cloud computing [65] represents a centralized infrastructure supplying remote server, storage, and database resources to support AIGC service lifecycle processes, including data collection, model training, fine-tuning, and inference. Cloud computing allows users to access AIGC services through the core network where these services are deployed, rather than building and maintaining physical infrastructure. Specifically, there are three primary delivery models in cloud computing: Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service (SaaS). In mobile AIGC networks, IaaS providers offer access to virtualized AIGC computing resources such as servers, storage, and databases [19]. Additionally, PaaS provides a platform for developing and deploying AIGC applications and services. Lastly, SaaS delivers applications and services over the internet, enabling users to access AIGC models directly through a web browser or mobile application. In summary, cloud computing in mobile AIGC networks allows developers and users to harness the benefits of AI while reducing costs and mitigating challenges associated with constructing and maintaining physical infrastructure,

playing a critical role in the development, deployment, and management of AIGC services.

*2) Edge Computing:* By providing computing and storage infrastructure at the edge of the core network [25], users can access AIGC services through radio access networks (RAN). Unlike the large-scale infrastructure of cloud computing, edge servers' limited resources often cannot support AIGC model training. However, edge servers can offer real-time fine-tuning and inference services that are less computationally and storage-intensive. By deploying edge computing at the network's periphery, users need not upload data through the core network to cloud servers to request AIGC services. Consequently, reduced service latency, improved data protection, increased reliability and decreased bandwidth consumption are benefits of AIGC services delivered via edge servers. Compared to exclusively delivering AIGC services through centralized cloud computing, location-aware AIGC services at the edge can significantly enhance user experience [66]. Furthermore, edge servers for local AIGC service delivery can be customized and personalized to meet user needs. Overall, edge computing enables users to access high-quality AIGC services with lower latency.

*3) Device-to-device Mobile Computing:* Device-to-device (D2D) mobile computing involves using mobile devices for the direct execution of AIGC services by users [14]. On one hand, mobile devices can directly execute AIGC models and perform local AIGC inference tasks. While running AIGC models on devices demands significant computational resources and consumes mobile device energy, it reduces AIGC service latency and protects user privacy. On the other hand, mobile devices can offload AIGC services to edge or cloud servers operating over wireless connections, providing a flexible scheme for delivering AIGC services. However, offloading AIGC services to edge or cloud servers for execution necessitates stable network connectivity and increases service latency. Lastly, model compression and quantization must be considered to minimize the resources required for execution on mobile devices, as AIGC models are often large-scale.

### D. Lessons Learned

*1) Cloud-Edge Collaborative Training and Fine-tuning for AIGC models:* To support AIGC services with required performance evaluated based on metrics discussed in Section III-A, cloud-edge collaborative pre-training, and fine-tuning are envisioned to be promising approaches. On the one hand, the servers in cloud computing can train AIGC models by using powerful computing and data resources. On the other hand, based on the large amount of user data in the edge network, the AIGC model can be fine-tuned to be more customized and personalized.

*2) Edge-Mobile Collaborative Inference for AIGC Services:* In a mobile AIGC network, the user's location and mobility change over time. Therefore, a large number of edge and mobile collaborations are required to complete the provision of AIGC inference services. Due to the different mobility of users, the AIGC services forwarded to the edge servers for processing are also dynamic. Therefore, dynamic resource

(a) Stable Diffusion      (b) DALLE-2      (c) Visual ChatGPT      (d) Point-E

Fig. 7: Generated images of different AIGC models, including Stable Diffusion (https://huggingface.co/spaces/stabilityai/stable-diffusion), DALLE-2 (https://labs.openai.com/), Visual ChatGPT (https://huggingface.co/spaces/microsoft/visual_chatgpt), Point-E (https://huggingface.co/spaces/openai/point-e), using the prompt "A photo of a green pumpkin".

allocation and task offloading decisions in mobile AIGC networks are one of the challenges in deploying mobile AIGC networks, which we discuss in Section VI.

## IV. How to Deploy AIGC at Mobile Edge Networks: Applications and Advantages of AIGC

This section introduces creative applications and advantages of AIGC services in the mobile edge network. Then, we provide four use cases of AIGC applications at mobile AIGC networks. Some examples of AIGC models are shown in Fig. 7. The applications elaborated in this section are summarized in Table II.

### A. Applications of Mobile AIGC Networks

*1) AI-generated Texts:* Recent advancements in Natural Language Generation (NLG) technology have led to AI-generated text that is nearly indistinguishable from human-written text [9]. The availability of powerful open-source AI-generated text models, along with their reduced computing power requirements, has facilitated widespread adoption, particularly in mobile networks. The development of lightweight NLG models that can operate on resource-constrained devices, such as smartphones and IoT devices, while maintaining high-performance levels, has made AI-generated text an essential service in mobile AIGC networks [34].

One example of such a model is ALBERT (A Lite BERT), designed to enhance the efficiency of BERT (Bidirectional Encoder Representations from Transformers) while reducing its computational and memory requirements [101]. ALBERT is pre-trained on a vast corpus of text data and uses factorized embedding parameterization, cross-layer parameter sharing, and sentence-order prediction tasks to optimize BERT's performance while minimizing computational and memory demands. ALBERT has achieved performance levels comparable to BERT on various natural language processing tasks, such as question answering and sentiment analysis [10]. Its lighter model design makes it more suitable for deployment on edge devices with limited resources.

MobileBERT is another model designed for deployment on mobile and edge devices with minimal resources [102]. This more compact variant of the BERT model is pre-trained on the same amount of data as BERT but features a more computationally efficient design with fewer parameters. Quantization is employed to reduce the model's weight accuracy, further decreasing its processing requirements. MobileBERT is a highly efficient model compatible with various devices, including smartphones and IoT devices, and can be used in multiple mobile applications, such as personal assistants, chatbots, and text-to-speech systems [34]. Additionally, it can be employed in small-footprint cross-modal applications, such as image captioning, video captioning, and voice recognition. These AI-generated text models offer significant advantages to mobile edge networks, enabling new applications and personalized user experiences in real time while preserving user privacy.

*2) AI-generated Audio:* AI-generated audio has gained prominence in mobile networks due to its potential to enhance user experience, and increase efficiency, security, personalization, cost-effectiveness, and accessibility [16]. For instance, AIGC-based speech synthesis and enhancement can improve call quality in mobile networks, while AIGC-based speech recognition and compression can optimize mobile networks by reducing the data required to transmit audio and automating tasks such as speech-to-text transcription. Voice biometrics powered by AI can bolster mobile network security by utilizing the user's voiceprint as a unique identifier for authentication [93]. AIGC-driven audio services, such as personalized music generation, can automate tasks and reduce network load, thereby cutting costs.

Audio Albert [41], a streamlined version of the BERT model adapted for self-supervised learning of audio representations, demonstrates competitive performance levels compared to other popular AI-generated audio models in various natural language processing tasks such as speech recognition, speaker identification, and music genre classification. In terms of latency, Audio Albert shows faster inference times than previous models, with a 20% reduction in average inference time on average, which can significantly improve response times in mobile edge networks. Additionally, Audio Albert's accuracy is comparable to BERT and achieves state-of-the-art results on several benchmarks. Furthermore, Audio Albert's model design is lighter than other models, making it suitable for

TABLE II: Summary of State-of-the-art AIGC models.

| Application | Models | Network Architectures | Datasets | Evaluation Metrics |
|---|---|---|---|---|
| Text Generation | GPT-3 [67], GPT-4, BERT [68], LaMDA [69], ChatGPT [10] | Transformer [70] | WebText, BookCorpus [71], Common Crawl | BLEU [72], ROUGE [73], Perplexity |
| Image Generation | StyleGAN [74], BigGANs [75], StyleGANXL [76], DVD-GAN [77], DALLE [6], DALLE2 [7], CLIP [5], VisualGPT [78], VAE [79], Energy-based GAN [62], Flow-based models [64], Imagen [80], diffusion probabilistic models [81], DDPM [82], DDIM [83] | GAN [84], VQ-VAE [85], Transformer [70] | ImageNet [86], CelebA [87], COCO [88] | FID [89], IS [90], LPIPS [91] |
| Music Generation | MuseNet [92], Jukedeck, WaveNet [93], AudioLM [94] | Transformer, RNN, CNN | MIDI Dataset, MAESTRO [95] | ABC-notation, Music IS |
| Video Generation | Diffusion models beat GANs [96], Video Diffusion Models [97], Dreamfusion [98] | DDPM, DDIM | Kinetics [99] | PSNR, SSIM |
| 3D Generation | NeRF [100] | MLP | Synthetic and real-world scenes | PSNR, SSIM, LPIPS |

deployment on edge devices with limited resources, improving computational efficiency while maintaining high-performance levels. Utilizing Audio Albert in mobile edge networks can provide several benefits, such as faster response times, reduced latency, and lower power consumption, making it a promising solution for AI-generated audio in mobile edge networks.

*3) AI-generated Images:* AI-generated images offer numerous applications in mobile networks, such as image enhancement, image compression, image recognition, and text-to-image generation [103]. Image enhancement can improve picture quality in low-light or noisy environments, while image compression decreases the data required to transmit images, enhancing overall efficiency. Various image recognition applications include object detection, facial recognition, and image search. Text-to-image generation enables the creation of images from textual descriptions for visual storytelling, advertising, and virtual reality/augmented reality (VR/AR) experiences [104]–[106].

Make-a-Scene, a novel text-to-image generation model proposed in [107], leverages human priors to generate realistic images based on textual descriptions. The model consists of a text encoder, an image generator, and a prior human module trained on human-annotated data to incorporate common sense knowledge. In mobile networks, this model can be trained on a large dataset of images and textual descriptions to swiftly generate images in response to user requests, such as creating visual representations of road maps. This approach complements the techniques employed in [108] for generating images with specific attributes.

Furthermore, the Semi-Parametric Neural Image Synthesis (SPADE) method introduced in [108] generates new images from existing images and their associated attributes using a neural network architecture. This method produces highly

realistic images conditioned on input attributes and can be employed for image-to-image translation, inpainting, and style transfer in mobile networks. The SPADE method shares similarities with the text-to-image generation approach in [107], where both techniques focus on generating high-quality, realistic images based on input data.

However, the development of AI-generated image technology also raises concerns around deep fake technology, which uses AI-based techniques to generate realistic photos, movies, or audio depicting nonexistent events or individuals, as discussed in [13]. Deep fakes can interfere with system performance and affect mobile user tasks, leading to ethical and legal concerns that require more study and legislation.

*4) AI-generated Videos:* AI-generated videos, like AI-generated images, can be utilized in mobile networks for various applications, such as video compression, enhancement, summarization, and synthesis [77]. AI-generated videos offer several advantages over AI-generated images in mobile networks. They provide a more immersive and engaging user experience by dynamically conveying more information [109]. Moreover, AI-generated videos can be tailored to specific characteristics, such as style, resolution, or frame rate, to improve user experience or create videos for specific purposes, such as advertising, entertainment, or educational content [97]. Furthermore, AI-generated videos can generate new content from existing videos or other types of data, such as images, text, or audio, offering new storytelling methods [97].

Various models can be employed to achieve AI-generated videos in mobile networks. One such model is Imagen Video, presented in [11], which is a text-conditioned video generation system based on a cascade of video diffusion models. Imagen Video generates high-definition videos from text input using a base video generation model and an interleaved sequence

of spatial and temporal video super-resolution models. The authors describe the process of scaling up the system as a high-definition text-to-video model, including design choices such as selecting fully-convolutional temporal and spatial super-resolution models at specific resolutions and opting for v-parameterization for diffusion models. They also apply progressive distillation with classifier-free guidance to video models for rapid, high-quality sampling [11], [97]. Imagen Video not only produces high-quality videos but also boasts a high level of controllability and world knowledge, enabling the generation of diverse videos and text animations in various artistic styles and with 3D object comprehension.

*5) AI-generated 3D:* AI-generated 3D content is becoming increasingly promising for various wireless mobile network applications, including AR and VR [110]. It also enhances network efficiency and reduces latency through optimal base station placement [111], [112]. Researchers have proposed several techniques for generating high-quality and diverse 3D content using deep learning (DL) models, some of which complement one another in terms of their applications and capabilities.

One such technique is the Latent-NeRF model, proposed in [113], which generates 3D shapes and textures from 2D images using the NeRF architecture. This model is highly versatile and can be used for various applications, such as 3D object reconstruction, 3D scene understanding, and 3D shape editing for wireless VR services. Another technique, the Latent Point Diffusion (LPD) model presented in [114], generates 3D shapes with fine-grained details while controlling the overall structure. LPD has been shown to create more diverse shapes than other state-of-the-art models, making it suitable for 3D shape synthesis, 3D shape completion, and 3D shape interpolation. The LPD model complements the Latent-NeRF approach by offering more diverse shapes and finer details.

Moreover, researchers in [115] proposed the Diffusion-SDF model, which generates 3D shapes from natural language descriptions. This model utilizes a combination of voxelized signed distance functions and diffusion-based generative models, producing high-quality 3D shapes with fine-grained details while controlling the overall structure. This technique accurately generates 3D shapes from natural language descriptions, making it useful for applications such as 3D shape synthesis, completion, and interpolation. It shares similarities with the Latent-NeRF and LPD models in terms of generating high-quality 3D content [116].

### B. Advantages of Mobile AIGC

We then discuss several advantages of generative AI in mobile networks.

*1) Efficiency:* Generative AI models offer several efficiency benefits in mobile networks. One of the primary advantages is automation. Generative AI models can automate the process of creating text, images, and other types of media, reducing the need for human labor and significantly boosting productivity [117]. The outputs of generative models can be generated quickly and with minimal human intervention. This

is particularly beneficial for tasks such as data augmentation in mobile networks, where a substantial amount of synthetic data is required to train ML models for applications like object recognition or network optimization. Moreover, generative AI models can be implemented at the edge of mobile networks [118], [119], allowing them to produce data locally on devices like smartphones and IoT sensors. This is especially advantageous for tasks that demand generating a large volume of data, such as image and video synthesis for AR applications. Local data production can reduce the amount of data transmitted over the mobile network, alleviating network congestion, and enhancing the system's responsiveness and efficiency [39]. This results in improved user experiences and reduced latency in mobile applications that rely on real-time data generation and processing.

*2) Reconfigurability:* The reconfigurability of AIGC in mobile networks is a significant advantage. By deploying AI models in mobile networks, AIGC can produce a vast array of content, including text, images, and audio, which can be seamlessly adjusted to suit evolving network demands and user preferences [120]. For instance, the ChatGPT model exemplifies AIGC's reconfigurability in providing multilingual support. It can be trained to understand and address user queries in numerous languages, facilitating seamless system adaptation for handling various linguistic contexts. This approach showcases how AIGC can cater to diverse user bases and adapt to global communication needs in mobile networks.

However, implementing multilingual support in AIGC models can pose several challenges and limitations, such as the need for large amounts of training data and the difficulty of maintaining consistency across different languages. AIGC models require large amounts of training data to learn multiple languages, which can be difficult to obtain for less commonly spoken languages. Additionally, maintaining consistency across different languages can be challenging, as each language has its own unique grammar and syntax. This can lead to errors or inaccuracies in translation, especially for more complex language structures. Finally, AIGC models may struggle with cultural nuances, metaphors, and idiomatic expressions that are specific to certain languages, which can result in misunderstandings or misinterpretations of user queries. To overcome the challenges of implementing multilingual support in AIGC models, future research could focus on developing more efficient and effective training methods that require less data while still producing accurate results. Additionally, ongoing efforts to improve natural language processing and machine translation algorithms could help improve consistency across different languages and reduce errors in translation. Another potential solution is to incorporate cultural and linguistic experts into the training process to help AIGC models better understand cultural nuances and expressions specific to different languages. Finally, exploring the use of transfer learning, where a model trained on one language is adapted to another language with less training data, could also be a promising direction for future research.

Additionally, AIGC can contribute to reconfigurability in mobile networks through the utilization of image and audio generative models. These models can be trained to generate

new visuals and auditory content based on specific parameters, such as user preferences or contextual information. As a result, the mobile system can be rapidly altered to produce novel materials on demand, eliminating the need for manual labor or supplementary resources. Another potential application of AIGC is the development of dynamic network architectures in mobile networks. These AI-enhanced designs can be effortlessly reconfigured to address shifting network demands, such as fluctuations in traffic patterns or the introduction of innovative services. For example, generative AI models, such as diffusion models, can be used to create optimal system incentive mechanisms according to the network environment, thereby improving the utility of participating users and enhancing overall network performance [121].

*3) Accuracy:* Employing generative AI models in mobile networks provides significant benefits in terms of accuracy, leading to more precise predictions and well-informed decision-making [96]. Enhanced accuracy in AI-generated content can substantially improve the overall user experience across various applications within the mobile network ecosystem. For example, AI-generated text can automate responses to mobile user inquiries, augmenting the efficiency and precision of mobile user support. This application not only reduces response times but also ensures accurate and contextually relevant information is provided to users, leading to better customer satisfaction and streamlined support services [39]. Similarly, AI-generated visuals and audio can be employed to elevate the quality and accuracy of network-provided content, encompassing domains such as advertising, entertainment, and accessibility services. By using generative AI models, tailored and engaging content can be produced, resulting in a more impactful and personalized user experience. In the context of mobile networks, this can mean generating high-quality images or videos adapted to various devices and network conditions, improving the user's perception of the provided services. By harnessing the power of generative AI models, mobile networks can offer more accurate and efficient services, ultimately fostering a superior user experience and enabling innovative solutions tailored to the diverse needs of mobile users.

*4) Scalability and Sustainability:* Utilizing AIGC in mobile networks offers significant scalability and sustainability benefits [96]. AIGC can produce a wide range of content, including text, images, and audio, enhancing mobile networks' overall scalability and sustainability in numerous ways. Specifically, AIGC facilitates scalability in mobile networks by reducing the reliance on human labor and resources. For instance, AIGC can generate automated responses to customer inquiries, alleviating the need for human customer support staff. This approach decreases the energy consumption associated with operating human-staffed contact centers and reduces the carbon footprint linked to human labor [21]. Furthermore, AIGC can promote sustainability in mobile networks by diminishing the demand for physical content storage. By generating new content on demand, AIGC minimizes the necessity to store and manage vast quantities of physical materials. This reduction leads to decreased energy usage and a smaller carbon footprint tied to maintaining physical storage infrastructure. Despite the challenges associated with AIGC models, such as large model sizes and complex training processes, leveraging edge servers in mobile networks can help mitigate these issues by adopting an "AIGC-as-a-Service" approach [118]. Users can interact with the system by submitting requests through their mobile devices and subsequently receiving computational results from edge servers. This strategy eliminates the necessity to deploy AIGC models on devices with constrained computing resources, optimizing overall efficiency and further improving scalability and sustainability within the mobile network infrastructure.

*5) Security and Privacy:* AIGC can offer potential security and privacy advantages by embedding sensitive information within AI-generated content. This approach can serve as a form of steganography, a technique that conceals data within other types of data, making it difficult for unauthorized parties to detect the hidden information. For instance, AI-generated images or audio can be used to encode confidential information in imperceptible ways. This technique can improve privacy in mobile networks, as sensitive data can be transmitted without being explicitly discernible. In addition, AI-generated content can be employed as a security measure, such as AI-generated audio for voice biometrics or AI-generated facial images for authentication purposes, adding an extra layer of security to mobile network services [21]. However, it is essential to be aware of potential security and privacy risks associated with AIGC, such as adversarial attacks on AI models or the misuse of AI-generated content for malicious purposes, like deepfakes [13]. To ensure the secure and privacy-preserving use of AIGC in mobile networks, robust security measures and encryption techniques must be in place, along with ongoing research to counter potential threats [122].

## V. Case Studies of AIGC in Mobile Network

In this section, we present several case studies for mobile AIGC networks. Specifically, we discuss the AIGC service provider (ASP) selection, generative AI-empowered traffic and driving simulation, AI-generated incentive mechanism, and blockchain-powered lifecycle management for AIGC.

### A. AIGC Service Provider Selection

The integration of AIGC models within wireless networks offers significant potential, as these state-of-the-art technologies have exhibited exceptional capabilities in generating a wide range of high-quality content. By harnessing the power of artificial intelligence, AIGC models can astutely analyze user inputs and produce tailored, contextually relevant content in real-time [96]. This stands to considerably enhance user experience and foster the creation of innovative applications across various domains, such as entertainment, education, and communication. Nonetheless, the deployment and application of these advanced models give rise to challenges, including extensive model sizes, complex training processes, and resource constraints. Consequently, deploying large-scale AI models on every network edge device poses considerable difficulties.

To address this challenge, the authors in [118] introduce the "AIGC-as-a-service" architecture. This approach entails ASPs deploying AI models on edge servers, which facilitates
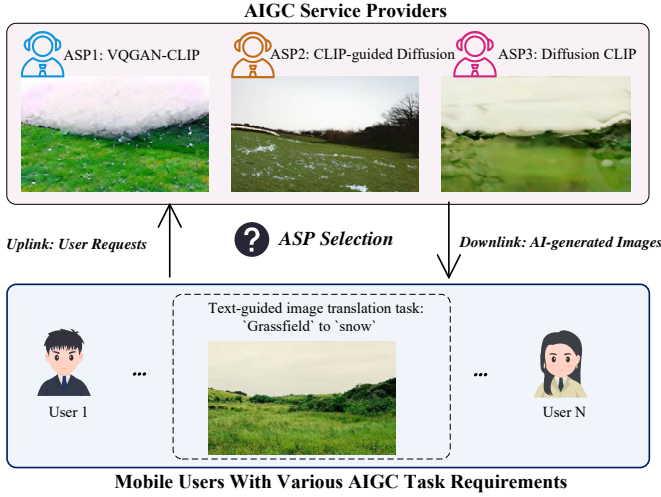
Fig. 8: The system model of AIGC service provider selection. Different ASPs performing user tasks can bring different results and different user utilities. Considering that different mobile users have different task requirements and different ASP's AI models have different capabilities and computation capacities, a proper ASP selection algorithm is needed to maximize the total utilities of network users.

the provision of instantaneous services to users via wireless networks, thereby ensuring a more convenient and adaptable experience. By enabling users to effortlessly access and engage with AIGC, the proposed solution minimizes latency and resource consumption. Consequently, edge-based AIGC-as-a-service holds the potential to transform the creation and delivery of AIGC across wireless networks.

However, one problem is that the effectiveness of ASP in meeting user needs displays significant variability due to a variety of factors. Certain ASPs may concentrate on generating specific content types, while others boast more extensive content generation capabilities. For instance, some providers may specialize in producing particular content categories, whereas others offer a wider range of content generation options. Moreover, several ASPs may have access to advanced computing and communication resources, empowering them to develop and deploy more sophisticated AIGC models within the mobile network. As depicted in Fig. 8, users uploading images and requirement texts to different ASPs encounter diverse results owing to the discrepancies in models employed. For example, a user attempting to add snow to grass in an image may experience varying outcomes depending on the ASP chosen.

With a large number of mobile users and increasing demand for accessing requests, it is crucial to analyze and select ASPs with the necessary capability, skill, and resources to offer high-quality AIGC services. This requires a rigorous selection process considering the provider's AIGC model capabilities and computation resources. By selecting a provider with the appropriate abilities and resources, organizations can ensure that they have effective AIGC services to increase the QoE for mobile users. Motivated by the aforementioned reasons, the authors in [118] examine the viability of large-scale

deployment of AIGC-as-a-Service in wireless edge networks. Specifically, in the ASP selection problem, which can be framed as a resource-constrained task assignment problem, the system consists of a series of sequential user tasks, a set of available ASPs, and the unique utility function for each ASP. The objective is to find an assignment of tasks to ASPs, such that the overall utility is maximized. Note that the utility of the task assigned to the ASP is a function of the required resource. Without loss of generality, the authors in [118] consider that is in the form of the diffusion step of the diffusion model, which is positively correlated to the energy cost. The reason is that each step of the diffusion model has energy consumption as it involves running a neural network to remove Gaussian noise. Finally, the total availability of resources for each ASP is taken into account to ensure that the resource constraints are satisfied.

In this formulation of AIGC service provisioning, the resource constraints are incorporated through the resource constraint, which specifies the limitations on the available resources. Note that failing to satisfy the resource constraint can result in the crash of ASP, causing the termination and restart of its running tasks.

Several baseline policies are used for comparison:

- **Random Allocation Policy.** This strategy distributes tasks to ASPs in a haphazard manner, without accounting for available resources, task duration, or any restrictions. The random allocation serves as a minimum benchmark for evaluating scheduling efficiency.
- **Round-Robin Policy.** The round-robin policy allocates tasks to ASPs sequentially in a repeated pattern. This approach can generate effective schedules when tasks are evenly distributed. However, its performance may be suboptimal when there are significant disparities among them.
- **Crash-Avoid Policy.** The crash-avoid policy prioritizes ASPs with greater available resources when assigning tasks. The goal is to prevent overburdening and maintain system stability.
- **Upper Bound Policy.** In this hypothetical scenario, the scheduler has complete knowledge of the utility each ASP offers to every user before task distribution. The omniscient allocation strategy sets an upper limit on the performance of user-centric services by allocating tasks to ASPs with the highest utility and avoiding system failures. However, this approach relies on prior information about the unknown utility function, which is unrealistic in practice.

The authors in [118] employed a Deep Reinforcement Learning (DRL) technique to optimize Application Service Provider (ASP) selection. In particular, they implemented the Soft Actor-Critic (SAC) method, which alternates between evaluating and improving the policy. Unlike traditional actor-critic frameworks, the SAC approach maximizes a balance between expected returns and entropy, allowing it to optimize both exploitation and exploration for efficient decision-making in dynamic ASP selection scenarios. To conduct the simulation, the authors consider 20 ASPs and 1000 edge users. Each ASP
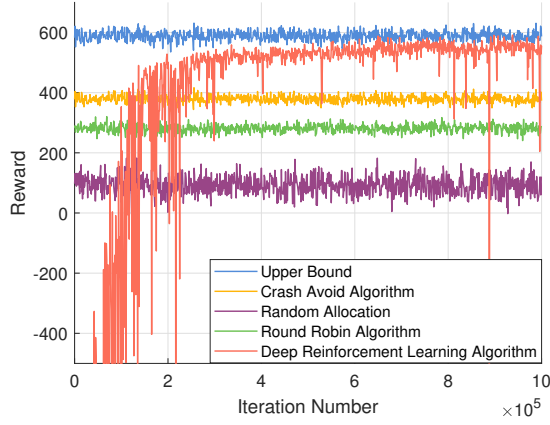
Fig. 9: The cumulative rewards under different ASP selection algorithms [118]. DRL-based algorithms can outperform multiple baseline policies, i.e., overloading-avoidance, random, and round-robin, and approximate the optimal policy.
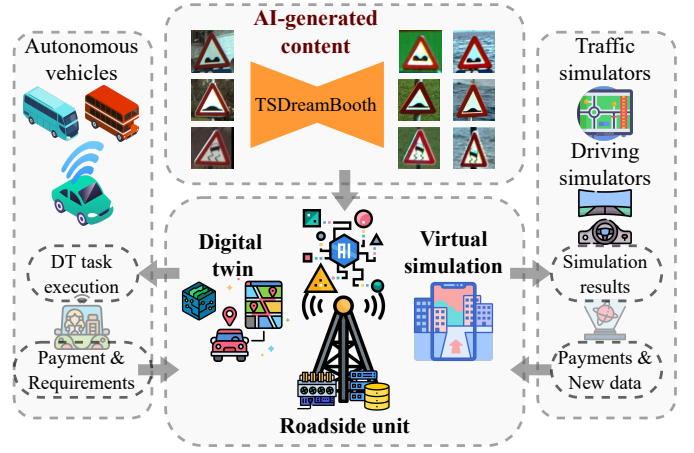


Fig. 10: Generative AI-empowered simulations for autonomous driving in vehicular Metaverse, which consists of AVs, virtual simulators, and roadside units.

offered AaaS with a maximum resource capacity, measured by total diffusion timesteps in a given time frame, varying randomly between 600 and 1,500. Each user submits multiple AIGC task requests to ASPs at varying times. These requests detailed the necessary AIGC resources in terms of diffusion timesteps, randomly set between 100 and 250. Task arrivals from users adhered to a Poisson distribution, with a rate of 0.288 requests per hour over a 288-hour duration, amounting to 1,000 tasks in total. As shown in Fig. 9, simulation results indicate that the proposed DRL-based algorithm outperforms three benchmark policies, i.e., overloading-avoidance, random, and round-robin, by producing higher-quality content for users and achieving fewer crashed tasks. ***Lesson Learned:*** The lesson learned from this study is that the proper selection of ASPs is crucial for maximizing the total utilities of network users and enhancing their experience. The authors in [118] introduced a DRL-based algorithm for ASP selection, which outperforms other baseline policies, such as overloading-avoidance, random, and round-robin. By leveraging the SAC approach, the algorithm strikes a balance between exploitation and exploration in decision-making for dynamic ASP selection scenarios. Consequently, this method can provide higher-quality content for users and lead to fewer crashed tasks, ultimately improving the quality of service in wireless edge networks. To further enhance research in the area of AIGC service provider selection, future studies could have:

- Investigate the integration of federated learning and distributed training methods to improve the efficiency of AIGC model updates and reduce the communication overhead among ASPs.
- Explore advanced DRL algorithms and meta-learning techniques to adaptively adjust the ASP selection strategy in response to changing network conditions and user requirements.
- Assess the impact of real-world constraints, such as network latency, data privacy, and security concerns, on the ASP selection process and devise strategies to address

these challenges.
- Develop multi-objective optimization techniques for ASP selection that consider additional factors, such as energy consumption, cost, and the trade-off between content quality and computational resources.

### B. Generative AI-empowered Traffic and Driving Simulation

In autonomous driving systems, traffic and driving simulation can affect the performance of connected autonomous vehicles (AVs). Existing simulation platforms are established based on historical road data and real-time traffic information. However, these data collection processes are difficult and costly, which hinders the development of fully automated transportation systems. Fortunately, generative AI-empowered simulations can largely reduce the cost of data collection and labeling by synthesizing traffic and driving data via generative AI models. Therefore, as illustrated in Fig. 10, the authors in [123] design a specialized generative AI model, namely TSDreambooth, for conditional traffic sign generation in the proposed vehicular mixed reality Metaverse architecture. In detail, TSDreambooth is a variation of stable diffusion [124] fine-tuned based on the Belgium traffic sign (BelgiumTS) dataset [125]. The performance of TSDreambooth is validated via the pre-trained traffic sign classification model as generative scores. In addition, the newly generated datasets are leveraged to improve the performance of original traffic sign classification models.

In the vehicular Metaverse, connected AVs, roadside units, and virtual simulators can develop simulation platforms in the virtual space collaboratively. Specifically, AVs maintain their representations in the virtual space via digital twin (DT) technologies. Therefore, AVs need to continuously generate multiple DT tasks and execute them to update the representations. To offload these DT tasks to roadside units for remote execution in real-time, AVs need to pay for the communication and computing resources of roadside units. Therefore, to provide fine-grained incentives for RSUs in executing DT tasks with heterogeneous resource demands and
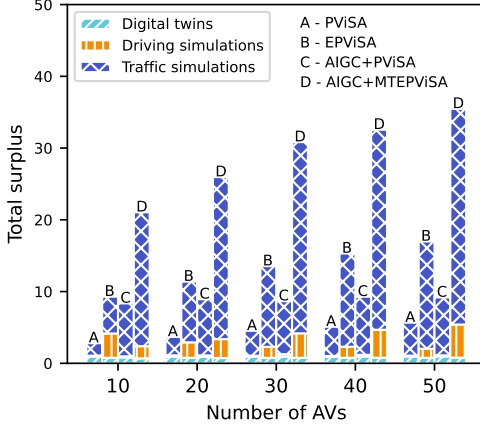
Fig. 11: Performance evaluation of the MTEPViSA under different sizes of the market.

various required deadlines, the authors in [123] propose a multi-task enhanced physical-virtual synchronization auction-based mechanism, namely MTEPViSA, to determine and price the resources of RSUs. There are two-stage of this mechanism the online submarket for provisioning DT services and the offline submarket for provisioning traffic and driving simulation services. In the online simulation submarket, the multi-task DT scoring rule is proposed to resolve the externalities from the offline submarket. In the meanwhile, the price scaling factor is leveraged to reduce the effect of asymmetric information among driving simulators and traffic simulators in the offline submarket. The simulation experiments are performed in a vehicular Metaverse system with 30 AVs, 30 virtual traffic simulators, 1 virtual driving simulator, and 1 RSU. The experimental results demonstrate that the proposed mechanism can improve 150% social surplus compared with other baseline mechanisms. Finally, they develop a simulation testbed of generative AI-empowered simulation systems in the vehicular Metaverse.

The vehicular mixed-reality (MR) Metaverse simulation environment was constructed employing a 3D model representing several city blocks within New York City. Geopipe, Inc. developed this model by leveraging artificial intelligence to generate a digital replica based on photographs taken throughout the city. The simulation encompasses an autonomous vehicle navigating a road, accompanied by strategically positioned highway advertisements. Eye-tracking data were gathered from human participants immersed in the simulation, utilizing the HMD Eyes addon provided by Pupil Labs. Subsequent to the simulation, participants completed a survey aimed at evaluating their subjective level of interest in each simulated scenario. As the experimental results shown in Fig. 11, According to the study, as the number of AVs continues to increase, the supply and demand mechanisms in the market are changing. Therefore, in order to improve market efficiency and total surplus, some mechanisms need to be adopted to coordinate supply and demand. We investigate the market mechanism and propose a mechanism based on AIGC technology to enhance market efficiency. Compared with the existing Physical-virtual Synchronization auction (PViSA) and Enhanced Physical-virtual Synchronization auction (EPViSA) mechanisms [126], [127], the AIGC-empowered mechanism can double the total surplus under different numbers of AVs.

***Lesson Learned:*** This case study on generative AI-empowered autonomous driving opens a new paradigm for the vehicular Metaverse, where data and resources can be utilized more efficiently. The authors demonstrate the potential of generative AI models in synthesizing traffic and driving data to reduce the cost of data collection and labeling. The proposed MTEPViSA mechanism also provides a solution to determine and price the resources of roadside units for remote execution of digital twin tasks, improving market efficiency and total surplus. However, there are still several open issues that need to be addressed in this field. Firstly, it is necessary to investigate the potential negative impacts of generative AI models in synthesizing traffic and driving data, such as biases and inaccuracies. Secondly, more research is needed to develop robust and trustworthy mechanisms for determining and pricing the resources of RSUs to ensure fair and efficient allocation of resources. Thirdly, the proposed mechanism needs to be tested and evaluated in more complex and varied scenarios to ensure its scalability and applicability in real-world situations.

### C. AI-Generated Incentive Mechanism

In this case study, we present the idea of using AI-generated optimization solutions with a focus on the use of diffusion models and their ability to optimize the utility function.

In today's world of advanced internet services, including the Metaverse, MR technology is essential for delivering captivating and immersive user experiences [128], [129]. Nevertheless, the restricted processing power of head-mounted displays (HMDs) used in MR environments poses a significant challenge to the implementation of these services. To tackle this problem, the researchers in [121] introduce an innovative information-sharing strategy that employs full-duplex device-to-device semantic communication [130]. This method enables users to circumvent computationally demanding and redundant processes, such as producing AIGC in-view images for all MR participants. By allowing a user to transmit generated content and semantic data derived from their view image to nearby users, these individuals can subsequently utilize the shared information to achieve spatial matching of computational outcomes within their own view images. In their work, the authors of [121] primarily concentrate on developing a contract theoretic incentive mechanism to promote semantic information exchange among users. Their goal is to create an optimal contract that, while adhering to the utility threshold constraints of the semantic information provider, simultaneously maximizes the utility of the semantic information recipient. Consequently, they devised a diffusion model-based AI-generated contract algorithm, as illustrated in Fig. 12.

Specifically, the researchers developed a cutting-edge algorithm for creating AI-generated incentive mechanisms, which tackle the challenge of utility maximization by devising optimal contract designs [121]. This approach is distinct from
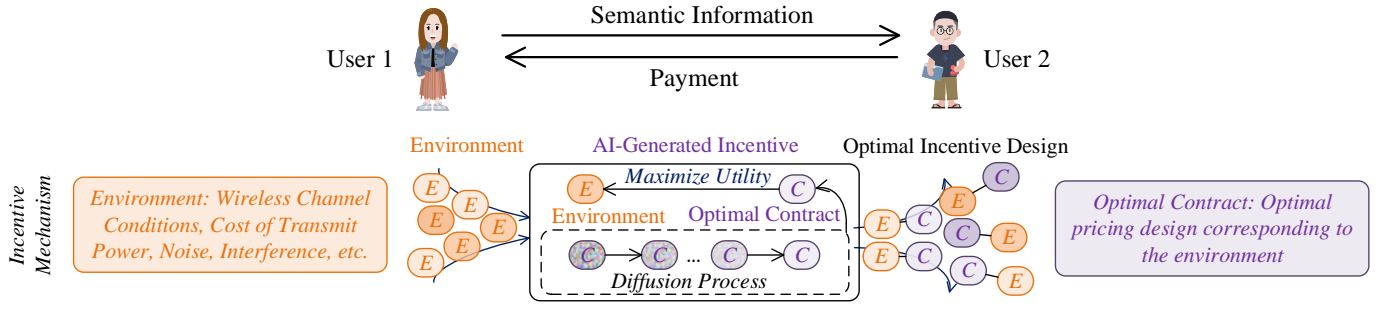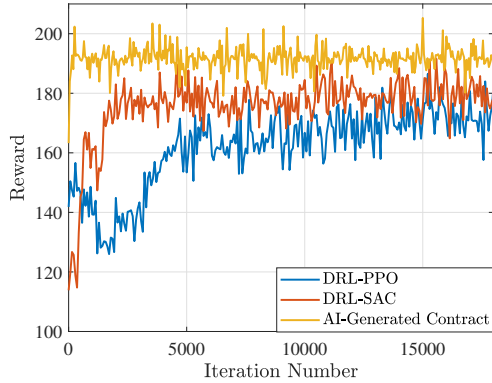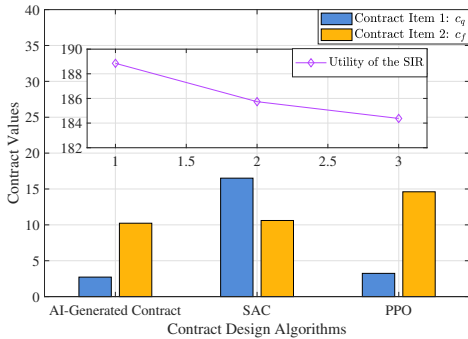
Fig. 12: System model of contract design in semantic information sharing network, and the AI-generated contract algorithm. The diffusion models generate different optimal contract designs under different environmental variables.



(a) Training process, with diffusion step $N = 10$ [121].



(b) The designed contracts.

Fig. 13: The effect of different incentive design schemes, e.g., PPO, SAC, and AI-generated contract [121].

with a value representing the expected total reward when an agent implements a particular contract design policy from the current state and adheres to it in the future. The optimal contract design policy is one that maximizes the system's predicted cumulative utility. The researchers then carry out an extensive comparison between their suggested AI-powered contract algorithm and two DRL algorithms, specifically SAC and PPO. As illustrated in the training process in [121] (see Fig. 13), PPO requires more iteration steps to achieve convergence, while SAC converges more quickly but with a lower final reward value in comparison to the AI-driven contract algorithm.

The enhanced performance of the suggested AI-driven contract algorithm can be ascribed to two main aspects:

- Improved sampling quality: By configuring the diffusion step to 10 and applying multiple refinement steps, the diffusion models generate higher quality samples, mitigating the influence of uncertainty and augmenting sampling precision [96].
- Enhanced long-term dependence processing capability: Unlike conventional neural network generation models that take into account only the current time step input, the diffusion model creates samples with additional time steps through numerous refinement iterations, thereby bolstering its long-term dependence processing capability [103].

As demonstrated in Fig. 13, the authors in [121] examine the optimal contract design capacities of the trained models. For a specific environmental state, the AI-driven contract algorithm provides a contract design that attains a utility value of 189.1, markedly outperforming SAC's 185.9 and PPO's 184.3. These results highlight the practical advantages of the proposed AI-based contract algorithm in contrast to traditional DRL techniques.

***Lesson Learned:*** The case study in this research highlights the potential of AI-generated optimization solutions, particularly diffusion models, for addressing complex utility maximization problems within incentive mechanism design. The authors in [121] present an innovative approach that employs full-duplex device-to-device semantic communication for information-sharing in mixed reality environments, overcoming the limitations of HMDs. The diffusion model-based AI-generated contract algorithm proposed in this study demon-

traditional neural network backpropagation algorithms or DRL methods, as it primarily focuses on enhancing contract design through iterative denoising of the initial distribution instead of optimizing model parameters. The policy for contract design is defined by the reverse process of a conditional diffusion model, linking environmental states to contract arrangements. The primary goal of this policy is to produce a deterministic contract design that maximizes the expected total reward over a series of time steps. To optimize system utility through contract design, the researchers in [121] create a contract quality network that associates an environment-contract pair

strates superior performance compared to traditional DRL algorithms, such as SAC and PPO. The superior performance of the AI-generated contract algorithm can be attributed to improved sampling quality and enhanced long-term dependence processing capability. This study underscores the effectiveness of employing AI-generated optimization solutions in complex, high-dimensional environments, particularly in the context of incentive mechanism design. Some promising directions for future research include:

- Expanding the application of diffusion models: Investigate the application of diffusion models in other domains, such as finance, healthcare, transportation, and logistics, where complex utility maximization problems often arise.
- Developing novel incentive mechanisms: Explore the development of new incentive mechanisms that combine AI-generated optimization solutions with other approaches, such as game theory or multi-agent reinforcement learning, to create even more effective incentive designs.
- Exploring the role of human-AI collaboration: Investigate how AI-generated optimization solutions can be combined with human decision-making to create hybrid incentive mechanisms that capitalize on the strengths of both human intuition and AI-driven optimization.

### D. Blockchain-Powered Lifecycle Management for AI-Generated Content Products

This case study delves into the application of a blockchain-based framework for managing the lifecycle of AIGC products within edge networks. The framework, proposed by the authors in [131], addresses concerns related to stakeholders, the blockchain platform, and on-chain mechanisms. We explore the roles and interactions of the stakeholders, discuss the blockchain platform's functions, and elaborate on the framework's on-chain mechanisms. Within edge networks, the AIGC product lifecycle encompasses four main stakeholders: content creators, Edge Service Providers (ESPs), end-users, and adversaries. The following describes their roles and interplay within the system:

- **Producers:** Initiate the AIGC product lifecycle by proposing prompts for ESPs to generate content. They retain ownership rights and can publish and sell the generated products.
- **ESPs:** Possess the resources to generate content for producers, charging fees based on the time and computing power used for the tasks.
- **Consumers:** View and potentially purchase AIGC products, participating in multiple trading transactions throughout the product lifecycle.
- **Attackers:** Seek to disrupt normal operations of AIGC products for profit through ownership tampering and plagiarism.

Considering the roles of these stakeholders, the blockchain platform fulfills two primary functions: providing a traceable and immutable ledger and supporting on-chain mechanisms. Transactions are recorded in the ledger and validated by full nodes using a consensus mechanism, ensuring security and traceability. ESPs act as full nodes, while producers and consumers serve as clients.

To address the concerns arising from stakeholder interactions, the framework employs three on-chain mechanisms [131]:

- **Proof-of-AIGC:** A mechanism that defends against plagiarism by registering AIGC products on the blockchain. It comprises two phases: proof generation and challenge.
- **Incentive Mechanism:** Safeguards the exchange of funds and AIGC ownership using Hashed Timelock Contracts (HTLCs).
- **Reputation-based ESP Selection:** Efficiently schedules AIGC generation tasks among ESPs based on their reputation scores.

The Proof-of-AIGC mechanism plays a vital role in maintaining the integrity of AIGC products. It encompasses two stages: proof generation and challenge. The objective of proof generation is to record AIGC products on the blockchain, while the challenge phase allows content creators to raise objections against any on-chain AIGC product they deem infringing upon their creations. If the challenge is successful, the duplicate product can be removed from the registry, thus protecting the original creator's intellectual property rights.

To further strengthen the security of the AIGC ecosystem, a pledge deposit is necessary to initiate a challenge, preventing arbitrary challenges that could burden the blockchain. This process comprises four steps: fetching the proofs, verifying the challenger's identity, measuring the similarity between the original product and the duplicate, and checking the results.

The AIGC economic system necessitates an incentive mechanism to motivate stakeholders and ensure legitimate exchanges of funds and ownership. The Incentive Mechanism rewards ESPs for maintaining the ledger and providing blockchain services. There are no transaction fees, and block generators follow a first-come-first-serve strategy. A two-way guarantee protocol using Hash Time Lock (HTL) is designed to build mutual trust and facilitate AIGC circulation during both the generation and trading phases.

The Proof-of-AIGC mechanism tackles issues like ownership manipulation and AIGC plagiarism, while the incentive mechanism ensures compliance with pre-established contracts. Furthermore, a reputation-based ESP selection accommodates ESP heterogeneity, which is crucial for efficient AIGC lifecycle management. Specifically, within the AIGC lifecycle management architecture, producers can concurrently interact with multiple heterogeneous ESPs, necessitating the identification of a trustworthy ESP for a specific task. Conventional approaches involve selecting the most familiar ESP to minimize potential risks, which may result in unbalanced workload distribution and increased service latency among ESPs. To address this challenge, a reputation-based ESP selection strategy is incorporated into the framework. This strategy ranks all accessible ESPs according to their reputation, which is computed using Multi-weight Subjective Logic (MWSL). The primary objectives are to assist producers in choosing the most reliable ESP, distribute the workload evenly across multiple ESPs, and motivate ESPs to accomplish tasks promptly and honestly, as a negative reputation impacts their earnings.
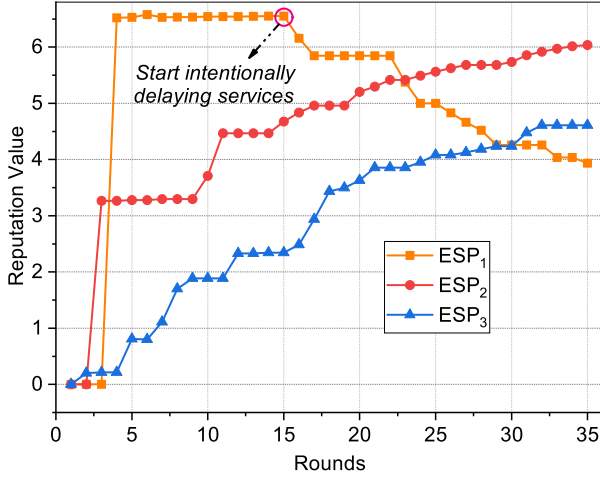
Fig. 14: The reputation trends of three ESPs (from the perspective of a random producer) [131].



Fig. 15: The total number of assigned tasks of three ESPs [131].

Producers identify suitable ESPs by computing the reputation of all potential ESPs, ranking them based on their current reputation, and allocating the AIGC generation task to the ESP with the highest standing. In MWSL, the concept of "opinion" serves as the fundamental element for reputation calculation. Local opinions represent the assessments of a specific producer who has directly interacted with the ESPs, while recommended opinions are derived from other producers who have also engaged with the ESPs. To mitigate the effect of subjectivity, an overall opinion is generated for each producer by averaging all the acquired recommended opinions. As producers possess varying degrees of familiarity with ESPs, the weight of their recommended opinions differs. Reputation is determined by combining a producer's local opinion with the overall opinion. The reputation scheme accomplishes its design objectives by quantifying the trustworthiness of ESPs, aiding producers in selecting the most dependable ESP, reducing service bottlenecks, and incentivizing ESPs to deliver high-quality AIGC services in order to maximize their profits.

A demonstration of the AIGC lifecycle management framework is conducted to verify the proposed reputation-based ESP selection approach [131]. The experimental setup comprises three ESPs and three producers, with the AIGC services facilitated by the Draw Things application. Several parameters are configured, and producers can employ the Softmax function to ascertain the probability of choosing each ESP. The reputation trends of the three ESPs are shown in Fig. 14, with ESP1 attaining the highest rank and remaining stable owing to its superior service quality. When ESP1 deliberately postpones AIGC services, its reputation declines sharply, while the reputations of ESP2 and ESP3 continue to rise. The proposed reputation strategy effectively measures the trustworthiness of ESPs, enabling producers to effortlessly discern the most reliable ESP and motivating ESPs to operate with integrity. The workload of ESPs under different ESP selection methods is also demonstrated in Fig. 15. Traditional methods lead to uneven 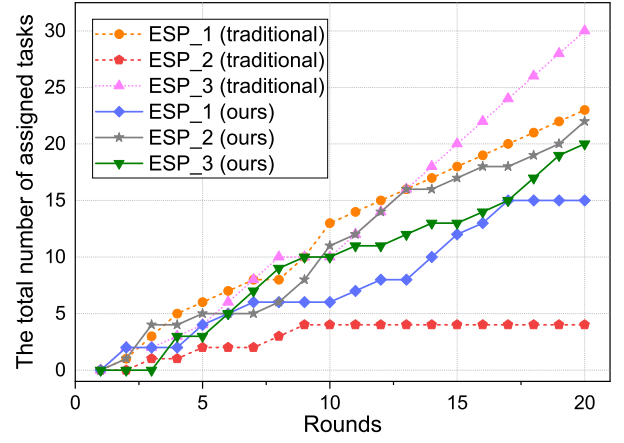workloads and extended service latencies. Conversely, the suggested reputation-based method efficiently balances the workload among ESPs, as producers can assess the trustworthiness of ESPs quantitatively without relying exclusively on their experiential judgment.

***Lesson Learned:*** The case study on blockchain-powered lifecycle management for AI-generated content products highlights the potential of a blockchain-based framework in addressing key concerns like stakeholder interactions, platform functionality, and on-chain mechanisms. The primary lessons learned emphasize the importance of defining clear stakeholder roles, implementing robust mechanisms such as Proof-of-AIGC and Incentive Mechanism to ensure system integrity, and employing a reputation-based ESP selection scheme to balance workload and encourage honest performance. These insights collectively contribute to the effective management of the AIGC product lifecycle within edge networks. Future research in blockchain-powered lifecycle management for AI-generated content products can explore several promising directions:

- Enhancing the efficiency and scalability of the blockchain platform to handle an increased number of transactions and support a growing AIGC ecosystem might be critical.
- Refining the reputation-based ESP selection scheme to account for more sophisticated factors, such as task complexity, completion time, and user feedback, could lead to more accurate and dynamic trustworthiness evaluations.
- Incorporating privacy-preserving techniques to protect sensitive data in AIGC products and user information without compromising the transparency and traceability of blockchain technology would be valuable.

## VI. IMPLEMENTATION CHALLENGES IN MOBILE AIGC NETWORKS

When providing AIGC services, a significant amount of computational and storage resources are required to run the AIGC model. These computation and storage-intensive services pose new challenges to existing mobile edge computing infrastructure. As discussed in Section III-C, a cloud-edge-mobile collaborative computing architecture can be implemented to provide AIGC services. However, several critical

implementation challenges must be addressed to improve resource utilization and the user experience.

## A. Edge Resource Allocation

AIGC service provisioning based on edge intelligence is computationally and communication-intensive for resource-constrained edge servers and mobile devices [132]. Specifically, AIGC users send service allocation requests to edge services. Upon receiving these AIGC requests, edge servers perform the AIGC tasks and deliver the output to users [133]. During this AIGC service provisioning interaction, model accuracy and resource consumption are the most common metrics. Consequently, significant efforts are being made to coordinate mobile devices and edge servers for deploying generative AI at mobile edge networks. As summarized in Table III, several Key Performance Indicators (KPIs) for edge resource allocation in AIGC networks are presented below. Here are several KPIs for edge resource allocation in AIGC networks.

- Model accuracy: In a resource-constrained edge computing network, a key issue when allocating edge resources is optimizing the accuracy of AI services while fully utilizing network resources. Besides objective image recognition and classification tasks, AI models are also based on the content's degree of personalization and adaptation. Thus, optimizing AIGC content networks may be more complex than traditional optimization since personalization and customization make evaluating model accuracy more unpredictable.
- Bandwidth utilization: While providing AIGC services, the edge server must maximize its channel utilization to ensure reliable service in a high-density edge network. To allocate its bandwidth resources more efficiently, the edge server must control channel access to reduce interference between user requests and maximize the quality of its AIGC service to attract more users.
- Edge resource consumption: Deploying AIGC services in edge networks requires computationally intensive AI training and inference tasks that consume substantial resources. Due to the heterogeneous nature of edge devices, edge services consume resources in generating appropriate AIGC while processing users' requests [141]. Deployment of AIGC services necessitates continuous iteration to meet actual user needs, as generation results of AIGC models are typically unstable. This constant AIGC service provisioning at edge servers leads to significant resource consumption.

Obtaining a balance between model accuracy and resource consumption can be challenging in resource-constrained edge computing networks. One potential strategy is to adjust the trade-off between model accuracy and resource consumption according to the needs of the users. For example, in some cases, a lower level of model accuracy may be acceptable if it results in faster response times or lower resource consumption. Another approach is to use transfer learning, which involves training an existing model on new data to improve accuracy while requiring fewer computational resources. Model
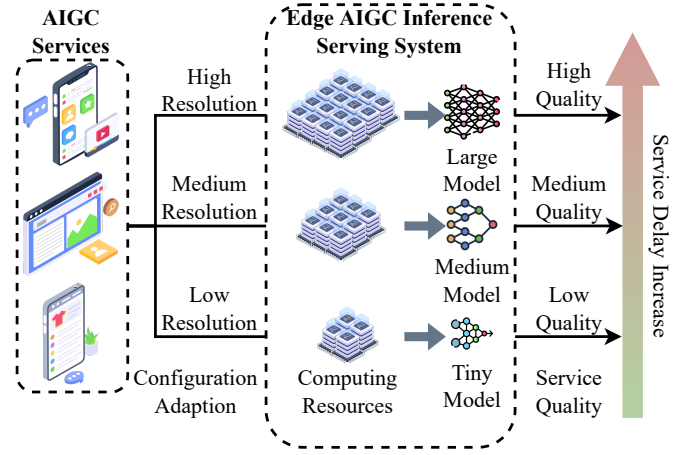


Fig. 16: Dynamic AIGC application configuration and AIGC model compression for serving AIGC services in mobile AIGC networks.

compression techniques can also be used to reduce the size of the AI model without significantly impacting accuracy. However, it is important to note that these techniques may not be applicable in all scenarios, as personalization and customization can make evaluating model accuracy more unpredictable. Deployment of AIGC services necessitates continuous iteration to meet actual user needs, as generation results of AIGC models are typically unstable. Due to the heterogeneous nature of edge devices, edge services consume resources in generating appropriate AIGC while processing users' requests. This constant AIGC service provisioning at edge servers leads to significant resource consumption.

To provide intelligent applications at mobile edge networks, considerable effort should focus on the relationship between model accuracy, networking, communication, and computation resources at the edge. Simultaneously, offering AIGC services is challenging due to the dynamic network environment and user requirements at mobile edge networks. The authors in [135] propose a threshold-based approach for reducing traffic at edge networks during collaborative learning. By considering computation resources, the authors in [134] examine the distributed ML problem under communication, computation, storage, and privacy constraints. Based on the theoretical results obtained from the distributed gradient descent convergence rate, they propose an adaptive control algorithm for distributed edge learning to balance the trade-off between local updates and global parameter aggregations. The experimental results demonstrate the effectiveness of their algorithm under various system settings and data distributions.

AIGC models often require frequent fine-tuning and retraining for newly generated data and dynamic requests in non-stationary mobile edge networks [142]. Due to limited storage resources at edge servers and the different customization demands of AIGC providers, the AIGC service placement problem is investigated in [136]. To minimize total time and energy consumption in edge AI systems, the AI service placement and resource allocation problem is formulated as an MINLP. In the optimization problem, AI service placement

TABLE III: Summary of scenarios, problems, benefits/challenges, and mathematical tools of edge resource allocation.

| Ref. | Scenarios | Performance Metrics/Decision Variables | Benefits/Challenges | Mathematical Tools |
|---|---|---|---|---|
| [134] | Adaptive control for distributed edge learning | Model loss/Steps of local updates, the total number of iterations | Provisioning AIGC services in resource-constrained edge environments | Control theory |
| [135] | Geo-distributed ML | Execution time/Selective barrier, mirror clock | Provisioning Localized AIGC services | Convergence analysis |
| [136] | AI service placement in mobile edge intelligence | Total time and energy consumption/Service placement decision, local CPU frequencies, uplink bandwidth, edge CPU frequency | Fully utilize scarce wireless spectrum and edge computing resources in provisioning AIGC services | ADMM |
| [137] | Joint model training and task inference | Energy consumption and execution latency/Model download decision and task splitting ratio | Integrated fine-tuning and inference for AIGC models with heterogeneous computing resources | ADMM |
| [138] | Serving edge DNN inference for multiple applications and multiple models | Inference accuracy, latency, resource cost/Application configuration, DNN model selection, and edge resources | Provision rich AIGC services for long-term utility maximization | Regularization-based online optimization |
| [139] | Multi-user collaborative DNN partitioning | Execution latency/Partitioning, computation resources | Providing insights for partitioning AIGC models under edge-mobile collaboration | Iterative alternating optimization |
| [140] | Hierarchical federated edge learning | Data convergence and revenue/Cluster selection and payment | Provisioning privacy-preserving AIGC services in edge networks | Evolutionary game and auction |

and channel allocation are discrete decision variables, while device and edge frequencies are continuous variables. However, solving this problem is not trivial, particularly in large-scale network environments. Thus, the authors propose an alternating direction method of multipliers (ADMM) to reduce the complexity of solving this problem. The experimental results demonstrate that this method achieves near-optimal system performance while the computational complexity grows linearly as the number of users increases. Moreover, when edge intelligence systems jointly consider AI model training and inference [137], the ADMM method can optimize edge resources. Additionally, the authors [138] explore how to serve multiple AI applications and AI models at the edge. They propose EdgeAdapter, as illustrated in Fig. 16, to balance the triple trade-off between inference accuracy, latency, and resource consumption. To provide inference services with long-term profit maximization, they first analyze the problem as an NP-hard problem and then solve it with a regularization-based online algorithm.

In mobile AIGC networks, an effective architecture for providing AIGC services is to partition a large AIGC model into multiple smaller models for local execution [28]. In [139], the authors consider a multi-user scenario with massive IoT [143] devices that cooperate to support an intelligent application collaboratively. Although partitioning large ML models and distributing smaller models to mobile devices for collaborative execution is feasible, the model distribution and result aggregation might incur extra latency during model training and inference. Additionally, the formulated optimization

problem is complex due to its numerous constraints and vast solution space. To address these issues, the authors propose an alternative iterative optimization to obtain solutions in polynomial time. Furthermore, AIGC services allow users to input their preferences into AIGC models. Therefore, to preserve user privacy among multiple users during collaborative model training and inference, the authors in [140] investigate the communication efficiency issues of decentralized edge intelligence enabled by FL. In the FL network, thousands of mobile devices participate in model training. However, selecting appropriate cluster heads for aggregating intermediate models can be challenging. Decentralized learning approaches can improve reliability while sacrificing some communication performance, unlike centralized learning with a global controller. A two-stage approach can be adopted in decentralized learning scenarios to improve the participation rate. In this approach, evolutionary game-based allocation can be used for cluster head selection, and DL-based auction effectively rewards model owners.

### B. Task and Computation Offloading

In general, executing AIGC models that generate creative and valuable content necessitates substantial computational resources, which is impractical for mobile devices with limited resources [21], [150]. Offering high-quality and low-latency AIGC services is challenging for mobile devices with low processing power and limited battery life. Fortunately, AIGC users can offload the tasks and computations of AIGC models

TABLE IV: Summary of scenarios, problems, benefits/challenges, and mathematical tools of task and computation offloading.

| Ref. | Scenarios | Performance Metrics/Decision variables | Benefits/Challenges | Mathematical Tools |
|---|---|---|---|---|
| [144] | Edge intelligence in IoT | Processing delay/Task offloading decisions | Offload AIGC tasks for improving inference accuracy | Optimization theory |
| [145] | Intelligent IoT applications | Processing time/Offloading decisions | Support on-demand changes for AIGC applications | Random forest regression |
| [28] | Collaborative intelligence between the cloud and mobile edge | Latency and energy consumption/DNN computation partitioning | Cloud and mobile edge collaborative intelligence for AIGC models | Greedy algorithm |
| [27] | Cloud-edge intelligence | Service response time/Task processing node | Reduce the average response time for multi-task parallel AIGC services | Genetic algorithm |
| [146] | Cost-driven offloading for DNN-Based applications | System costs/Number of layers | Minimize costs of AIGC services in a cloud-edge-end collaborative environment | Genetic algorithm based on particle swarm optimization |
| [147] | Industrial edge intelligence | A weighted sum of task execution time and energy consumption/Task assignment | Multi-objective optimization of large-scale AIGC tasks with multiple connected devices | Generative coding evolutionary algorithm |
| [148] | Computation offloading for ML web apps | Inference time/Pre-sending decisions | Reduce execution overheads of AIGC tasks with pre-sending snapshots | Hill climbing algorithm |
| [149] | Cooperative edge intelligence | Quality of experience/Offloading decisions | Enhance vertical-horizontal cooperation in multi-user AIGC co-inference scenarios | Federated multi-agent reinforcement learning |

over the RAN to edge servers located in proximity to the users. This alleviates the computational burden on mobile devices.

As listed in Table IV, several KPIs are specifically relevant to computation offloading in mobile AIGC networks:

- Service latency: Service latency refers to the delay associated with data input and retrieval as well as the model inference computations that users perform to generate AIGC. By offloading AIGC tasks from mobile devices, such as fine-tuning and inference, to edge servers for execution, the total latency in mobile AIGC networks can be reduced. Unlike local execution of the AIGC model, offloading AI tasks to the edge server for execution introduces additional latency when transmitting personalized instructions and downloading AIGC content.

- Reliability: Reliability evaluates users' success rate in obtaining personalized data accurately. On the one hand, when connecting to the edge server, users may experience difficulty uploading the requested data to edge servers or downloading the results from servers due to dynamic channel conditions and wireless network instability. On the other hand, the content generated by the AIGC model may not fully meet the needs of AIGC users in terms of personalization and customization features. Unsuccessful content reception and invalid content affect the AIGC network's reliability.

When implementing cloud-edge collaborative training and fine-tuning for AIGC models, it is important to consider specific algorithms or techniques that enable effective collaboration between cloud and edge servers [151], [152]. For example, federated learning and distributed training approaches can
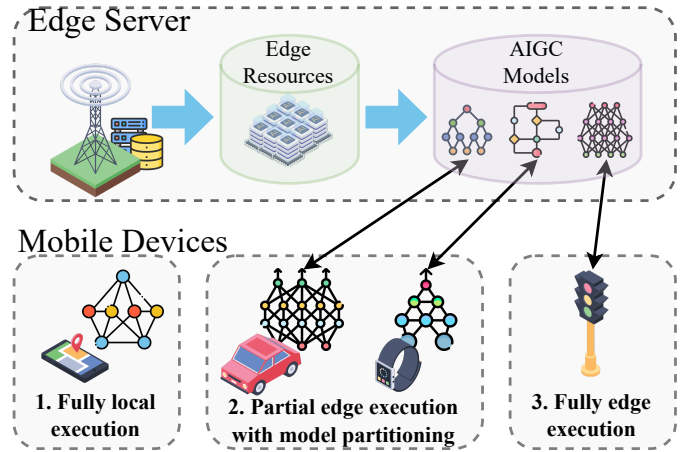


Fig. 17: Model partitioning in mobile AIGC networks. The AIGC models of mobile devices can be split and full or partial of them can be offloaded to edge servers for remote execution.

facilitate the collaboration process by allowing edge servers to train models locally and then send the updated weights to the cloud server for aggregation. The division of responsibilities between cloud and edge servers can also greatly affect the overall efficiency and performance of the AIGC models. Therefore, it is crucial to discuss and implement appropriate schemes for determining which tasks are offloaded to the edge servers and which are performed on the cloud server. To provide AIGC services in edge intelligence-empowered IoT, offloading ML tasks to edge servers for remote execution

is a promising approach for computation-intensive AI model inference. For instance, in Fig 17, multiple lightweight ML models can be loaded into IoT devices, while large-scale ML models can be installed and executed on edge servers [25]. Heterogeneous AIGC models can be deployed on mobile devices and edge servers according to their resource demands and service requirements [153]. However, the multiple attributes of ML tasks, such as accuracy, inference latency, and reliability, render the offloading problem of AIGC highly complex. Therefore, the authors in [144] propose an ML task offloading scheme to minimize task execution latency while guaranteeing inference accuracy. Considering error inference leading to extra delays in task processing, they initially model the inference process as M/M/1 queues, which are also applicable to the AIGC service process. Furthermore, the optimization problem of ML task execution is formulated as a Mixed-Integer Nonlinear Programming (MINLP) to minimize provisioning delay, which can be adopted in the inference process of AIGC services. To extend the deterministic environment in [144] into a more general environment, the authors in [145] first propose an adaptive translation mechanism to automatically and dynamically offload intelligent IoT applications. Then, they make predictive offloading decisions using a random forest regression model. Their experiments demonstrate that the proposed framework reduces response times for complex applications by half. Such ML methods can also be used to analyze AIGC network traffic to improve service delivery efficiency and reliability.

The success of edge-mobile collaboration for AIGC services is dependent on several factors, including the type of service, user characteristics, computational resources, and network conditions [3]. For instance, a real-time AIGC service may have different latency requirements compared to an offline service. Similarly, the required computational resources may vary depending on the model's complexity [154]. Additionally, the user profile, including location and device type, may affect the selection of edge servers for task offloading. Furthermore, network conditions such as bandwidth and packet loss rate can impact the reliability and latency of the service. Therefore, it is necessary to implement effective resource allocation and task offloading schemes to ensure high-quality and low-latency AIGC services in dynamic and diverse environments. Cloud-edge collaborative intelligence enables local tasks to be offloaded to edge and cloud servers. AIGC can benefit from cloud-edge intelligence, as edge servers can provide low-latency AIGC services while cloud servers can offer high-quality AIGC services. The authors in [28] develop a scheme called Neurosurgeon to select the optimal partitioning point based on model architectures, hardware platforms, network conditions, and load information at the servers to automatically partition the computation of tensors of DNNs between cloud and edge servers. Furthermore, the authors in [155] find that the layered approach can reduce the number of messages transmitted between devices by up to 97% while only decreasing the accuracy of models by a mere 3%. However, multiple AIGC services should be considered in cloud-edge collaborative intelligence that differs in types (e.g., text, images, and videos) and their diverse quality of service

(QoS) requirements. In multi-task parallel scheduling [27], the genetic algorithm can also be used to make real-time model partitioning decisions. The authors in [146] propose a cost-driven strategy for AI application offloading through a self-adaptive genetic algorithm based on particle swarm optimization.

In industrial edge intelligence, where edge intelligence is embedded in the industrial IoT [147], offloading computation tasks to edge servers is an efficient solution for self-organizing, autonomous decision-making, and rapid response throughout the manufacturing lifecycle, which is similarly required by mobile AIGC networks. Therefore, efficiently solving task assignment problems is crucial for effective AIGC model inference. However, the coexistence of multiple tasks among devices makes system response slow for various tasks. For example, text-based and image-based AIGC may coexist on the same edge device. As one solution, in [147], the authors propose a coding group evolution algorithm to solve large-scale task assignment problems, where tasks span the entire lifecycle of various products, including real-time monitoring, complex control, product structure computation, multidisciplinary cooperation optimization, and production process computation. Likewise, the AIGC lifecycle includes data collection, labeling, model training and optimization, and inference. Furthermore, a simple grouping strategy is introduced to parallel partition the solution space and accelerate the evolutionary optimization process. In contrast to VM-level adaptation to specific edge servers [156], the authors propose application-level adaptation for generic servers. The lighter adaptation framework in [148] further improves transmission time and user data privacy performance, including offloading and data/code recovery to generic edge servers.

Ensuring dependable task offloading is crucial in providing superior AIGC services with minimal latency in edge computing. For instance, data transmission redundancy can enhance dependability by transmitting data via multiple pathways to mitigate network congestion or failures. By incorporating these techniques, task offloading dependability in edge computing can be enhanced, thereby leading to more efficient and effective AIGC services. Most intelligent computing offloading solutions converge slowly, consume significant resources, and raise user privacy concerns [157], [158]. The situation is similar when leveraging learning-based approaches to make AIGC service offloading decisions. Consequently, the authors enhance multi-user QoE [159] for cooperative edge intelligence in [149] with federated multi-agent reinforcement learning. They formulate the cooperative offloading problem as a Markov Decision Process (MDP). The state is composed of current tasks, local loads, and edge loads. Learning agents select task processing positions to maximize multi-user QoE, which simultaneously considers service latency, energy consumption, task drop rate, and privacy protection. Similarly, AIGC service provisioning systems can easily adopt the proposed solution for maximizing QoE in AIGC services.

### C. Edge Caching

Edge caching is the delivery of low-latency content and computing services using the storage capacity of edge base

TABLE V: Summary of scenarios, problems, performance metrics, and mathematical tools for edge caching in AIGC networks.

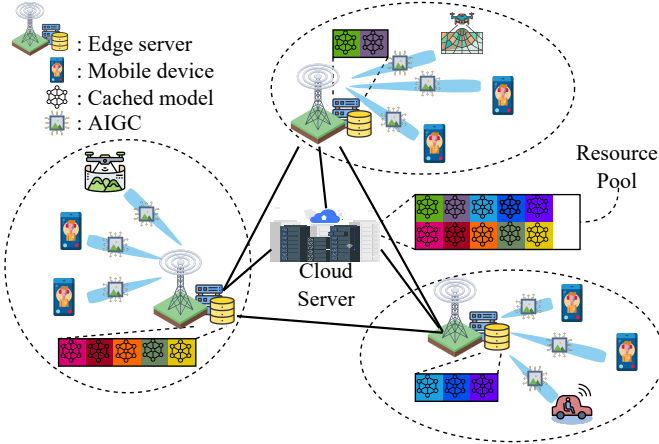| Ref. | Scenarios | Performance Metrics/Decision Variables | Benefits/Challenges | Mathematical Tools |
|---|---|---|---|---|
| [160] | DL Model caching at the edge | Runtime memory consumption and loading time/Model preload policy | Manage and utilize GPU memories of edge servers for caching AIGC models | Cache replacement algorithms |
| [161] | Caching many models at the edge | Model load and execution latency and monetary cost /Caching eviction policy | Improve scalability of mobile AIGC networks via model-level caching deployment and replacement | Model utility calculation |
| [162] | Cache for mobile deep vision applications | Latency, accuracy loss, energy saving/Caching policy, user selection, transmit power, bandwidth ratio | Caching for users' requests for multimodal AIGC services | Greedy algorithm |
| [163] | Cache for functions in serverless computing | Execution time, cold start proportion/Function keep-alive policy | Keep AIGC models alive and warm for in-contextual inference | Greedy-dual based approach |
| [164] | Knowledge caching for federated learning | Transmission latency and energy consumption/Caching policy, user selection, transmit power, bandwidth ratio | Privacy-preserving model caching via knowledge of AIGC requests | Optimization theory |



Fig. 18: An overview of edge caching in mobile AIGC networks. By caching the AIGC model on the edge servers, the latency of AIGC services can be reduced and the network congestion in the core network can be reduced.

stations and mobile devices [165]. As illustrated in Fig. 18, in mobile AIGC networks, users can request AIGC services without accessing cloud data centers by caching AIGC models in edge servers and mobile devices. Unlike the cache in traditional content distribution networks, the AIGC model cache also requires computing resources to support its execution. Additionally, the AIGC model needs to gather user historical requests and profiles in context to provide personalized services during the AIGC service process. As shown in Table V, here are several key performance indicators (KPIs) for edge caching in AIGC networks:

- Model access delay: Model access latency is an important indicator of AIGC service quality. The latency is lowest when the AIGC model is cached in the mobile device [166]. The model access latency must also be calculated considering the delay in the wireless communication network when the edge server provides the AIGC model. Finally, the core network latency must be considered when the cloud provides the AIGC service.
- Backhaul traffic load: The load on the backhaul traffic is significantly reduced, as the requests and results of AIGC services do not need to go through the core network when the AIGC model is cached in the mobile edge network.
- Model hit rate: Similar to content hit rate, the model hit rate is an important metric for AIGC models in the edge cache. It can be used for future model exits and loading during model replacement.

As there is sufficient infrastructure and resources in the cloud computing infrastructure, the AIGC model can be fully loaded into the GPU memory for real-time service requests. In contrast, the proposed EdgeServe in [160] keeps models in main memory or GPU memory so that they can be effectively managed or used at the edge. Similar to traditional CDNs, the authors use model execution caches at edge servers to provide immediate AI delivery. In detail, there are mainly three challenges in AIGC model caching:

- Constraint-memory edge servers: Compared to the resource-rich cloud, the resources of servers in the edge network, such as GPU memory, are limited [167]. Therefore, caching all AIGC models on one edge server is infeasible.
- Model-missing cost: When the mobile device user requests AIGC, the corresponding model is missed if the AIGC model used to generate the AIGC is not cached in the current edge server [161]. In contrast to the instantly available AIGC service, if the AIGC model is missing,

the edge server needs to send a model request to the cloud server and download the model, which causes additional overhead in terms of bandwidth and latency.

- Functionally equivalent models: The number of AIGC models is large and increases depending on the number of detailed tasks [168]. Meanwhile, AI models have similar functions in different applications, i.e., functionally equivalent. For example, for image recognition tasks, a large number of models with different architectures are proposed to recognize features in images, which have different model architectures and computation requirements.

To address these challenges, the authors in [160] formulate the problem of edge modeling as determining which DL models should be preloaded into memory and which should be discarded when the memory is full while satisfying the requirements of inferential response times. Fortunately, this edge model caching problem can be solved using existing cache replacement policies for edge content caching. The accuracies and computation complexities of DL models make this optimization problem more complicated than conventional edge caching problems. Similarly, for resource-constrained edge servers, the AIGC model can be dynamically deployed and replaced. However, an effective caching algorithm for loading and unloading the AIGC models to maximize the hit rate has not yet been investigated.

As the capabilities of AI services continue to grow and diversify, multiple models need to be deployed simultaneously at the edge to achieve various tasks, including classification, recognition, text/image/video generation [169]. Especially in mobile AIGC networks, multiple base models need to work together to generate a large amount of multimodal synthetic data. Many models play a synergistic role in the AIGC services at the edge of the network, while the support of multiple models also poses a challenge to the limited GPU memory of the edge servers. Therefore, the authors in [161] propose a model-level caching system with an eviction policy according to model characteristics and workloads. The model eviction policy is based on model utility calculation from cache miss penalty and the number of requests. This model-aware caching approach introduces a new direction for providing AIGC services at mobile edge networks with heterogeneous requests. Experimental results show that compared to the non-penalty-aware eviction policy, the model load delay can be reduced by 1/3. This eviction policy can also be adopted in the problem of which unpopular AIGC models should be unloaded.

At mobile AIGC networks, not only the AIGC model needs to be cached, but also the AIGC requests and results can be cached to reduce the latency of service requests in AIGC networks. To this end, the authors devise a principled cache design to accelerate the execution of CNN models by exploiting the temporal locality of video for continuous vision tasks to support mobile vision applications [170]. The authors in [162] propose a principled cache scheme, named DeepCache, to retrieve reusable results and reuse them within a fine-grained CNN by exploiting the temporal locality of the mobile video stream. In DeepCache, mobile devices do not need to offload any data to the cloud and can support the most popular models. Additionally, without requiring developers to

retrain models or tune parameters, DeepCache caches inference results for unmodified CNN models. Overall, DeepCache can reduce energy consumption by caching content to reduce model inference latency while sacrificing a small fraction of model accuracy.

In serverless computing for edge intelligence, mobile devices can call functions of AIGC services at edge servers, which is more resource-efficient compared to container and virtual machine (VM)-based AIGC services. Nevertheless, such functions suffer from the cold-start problem of initializing their code and data dependencies at edge servers. Although the execution time of each function is usually short, initialization, i.e., fetching and installing prerequisite libraries and dependencies before execution, is time-consuming [171]. Fortunately, the authors in [163] show that the caching-based keep-alive policy can be used to address the cold-start problem by demonstrating that the keep-alive function is equivalent to caching. Finally, to balance the trade-off between server memory utilization and cold-start overhead, a greedy dual-based caching algorithm is proposed.

Frequently, a large-scale AIGC model can be partitioned into multiple computing functions that can be efficiently managed and accessed during training, fine-tuning, and inference. FL models can be cached on edge servers to facilitate user access to instances and updates, thus addressing user privacy concerns [172], [173]. For example, the authors in [164] propose a knowledge cache scheme for FL in which participants can simultaneously minimize training delay and training loss according to their preference. Their insight is that there are two stimulations for caching knowledge for FL [174]: i) training data sufficiency and ii) connectivity stability. Experimental results show that the proposed preference-driven caching policy, based on the preferences (i.e., demands or desires for global models) of participants in FL, can outperform the random policy when user preferences are intense. Therefore, preference-based AIGC model caching should be extensively investigated for providing personalized and customized AIGC services at edge servers.

### D. Mobility Management

Mobile edge intelligence for the Internet of Vehicles and Unmanned Aerial Vehicle (UAV) networks relies on effective mobility management solutions [182]–[185] to provide mobile AIGC services. Furthermore, UAV-based AIGC service distribution offers advantages such as ease of deployment, flexibility, and extensive coverage for enhanced edge intelligence [186], [187]. Specifically, UAVs, with their line-of-sight communication links, can extend the reach of edge intelligence [188]. For example, flexible UAVs equipped with AIGC servers enable users to access AIGC services with ultra-low latency and high reliability, especially when fixed-edge servers are often overloaded in hotspot areas or expensive to deploy in remote areas, as illustrated in Fig. 19. In addition, UAV-enabled edge intelligence can be utilized to implement mobile AIGC content and service delivery.

As summarized in Table VI, here are several KPIs for mobility management in AIGC networks:

TABLE VI: Summary of scenarios, problems, benefits/challenges, and mathematical tools for mobility management.

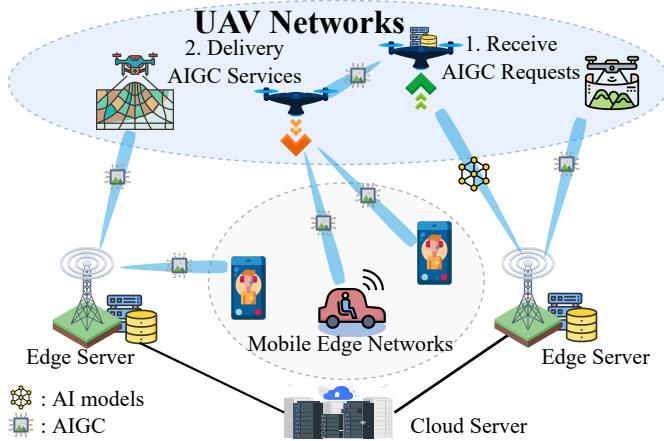| Ref. | Scenarios | Performance Metrics/Problems | Benefits/Challenges | Mathematical Tools |
|------|-----------|------------------------------|---------------------|--------------------|
| [175] | Jointing vehicle-edge deep neural network inference | Latency, failure rate/CPU frequency | Robust AIGC service provisioning via layer-level offloading | Chemical reaction optimization |
| [176] | Vehicular edge intelligence | Weighted average completion time and task acceptance ratio/Task dispatching policy | Provisioning AIGC service in multi-vehicle environments with motion prediction | Greedy algorithm |
| [177] | Mobility-enhanced edge intelligence | Task completion ratio and model accuracy/Offloading redundancy, task assignment, beam selection | Sustainable AIGC service provisioning with mobility management | Federated learning |
| [178] | Edge intelligence-assisted IoV | Average delay and energy consumption/Transmission decision, task offloading decision, bandwidth, and computation resource allocation | Flexible network model selection for AIGC services for balancing the tradeoff adaptively | Quantum-inspired reinforcement learning |
| [179] | Cooperative edge intelligence in IoV | Average delay and energy consumption/Trajectory prediction accuracy | Optimize AIGC service with spatial and temporal correlations of users' requests | Hybrid stacked autoencoder learning |
| [180] | UAVs as an intelligent service | Model accuracy and energy consumption/Number of local iterations | Provision AIGC services via a network of UAVs | Greedy algorithm |
| [181] | Knowledge distillation-empowered edge intelligence | Accuracy and inference delay/Size of model parameters | Visual information-aided AIGC model deployment and inference scheduling | Knowledge distillation |



Fig. 19: An overview of mobility management in mobile AIGC networks. The coverage of the mobile AIGC network will be significantly enhanced by UAV processing the user's server request and providing AIGC services.

- Task accomplishment ratio: The provisioning of AIGC services at mobile edge networks must consider the dynamic nature of users. As a result, services must be completed before users leave the base station. To measure the effectiveness of mobility management in AIGC networks, the task completion rate can be used.
- Coverage enhancement: Vehicles and UAVs can serve as reconfigurable base stations to enhance the coverage of mobile AIGC networks [189], providing AIGC models

and content to users anywhere and anytime.

In vehicular networks, intelligent applications, such as AIGC-empowered navigation systems, are reshaping existing transportation systems. In [175], the authors propose a joint vehicle-edge inference framework to optimize energy consumption while reducing the execution latency of DNNs. In detail, vehicles and edge servers determine an optimal partition point for DNNs and dynamically allocate resources for DNN execution. They propose a chemical reaction optimization-based algorithm to accelerate convergence when solving the resource allocation problem. This framework offers insights for implementing mobile AIGC networks, where vehicles can collaborate with base stations to provide real-time AIGC services based on DNNs during their movement.

AIGC applications require sufficient processing and memory resources to perform extensive AIGC services [190]–[193]. However, resource-constrained vehicles cannot meet the QoS requirements of the tasks. The authors in [176] propose a distributed scheduling framework that develops a priority-driven transmission scheduling policy to address the dynamic network topologies of vehicle networks and promote vehicle edge intelligence. To meet the various QoS requirements of intelligent tasks, large-volume tasks can be partitioned and sequentially uploaded. Additionally, the impact of vehicle motion on task completion time and edge server load balancing can be independently handled by intelligent task processing requests. The effectiveness of the proposed framework is demonstrated in single-vehicle and multi-vehicle environments through simulation and deployment experiments. To facilitate smart and green vehicle networks [177], the real-

time accuracy of AI tasks, such as AIGC model inference, can be monitored through on-demand model training using infrastructure vehicles and opportunity vehicles.

The heterogeneous communication and computation requirements of AIGC services in highly dynamic, time-varying Internet of Vehicles (IoV) warrant further investigation [194], [195]. To dynamically make transmission and offload decisions, the authors in [178] formulate a Markov decision process for time-varying environments in their joint communication and computation resource allocation strategy. Finally, they develop a quantum-inspired reinforcement learning algorithm, in which quantum mechanisms can enhance learning convergence and performance. The authors in [179] propose a stacked autoencoder to capture spatial and temporal correlations to combine road traffic management and data network traffic management. To reduce vehicle energy consumption and learning delay, the proposed learning model can minimize the required signal traffic and prediction errors. Consequently, the accuracy of AIGC services based on autoencoder techniques can be improved through this management framework.

With UAV-enhanced edge intelligence, UAVs can serve as aerial wireless base stations, edge computing servers, and edge caching providers in mobile AIGC networks. To demonstrate the performance of UAV-enhanced edge intelligence while preserving user privacy at mobile edge networks, the authors in [180] use UAV-enabled FL as a use case. Moreover, the authors suggest that flexible switching between compute and cache services using adaptive scheduling UAVs is a topic for future research. Therefore, flexible AIGC service provisioning and UAV-based AIGC delivery are essential for satisfying real-time service requirements and reliable generation. In this regard, the authors in [181] propose a visually assisted positioning solution for UAV-based AIGC delivery services where GPS signals are weak or unstable. Specifically, knowledge distillation is leveraged to accelerate inference speed and reduce resource consumption while ensuring satisfactory model accuracy.

### E. Incentive Mechanism

As suitable incentive mechanisms are designed, more edge nodes participate in and contribute to the AIGC services [131], [200], [201]. This increases the computational capacity of the system. In addition, the nodes are motivated to earn rewards by providing high-quality services. Thus, the overall quality of AIGC services is improved. Finally, nodes are encouraged to engage in secure operations without security concerns by recording resource transactions through the blockchain.

As listed in Table VII, here are several KPIs for incentive mechanisms in AIGC networks:

- Social welfare: AIGC's social welfare is the sum of the value of AIGC's services to the participants of the current network. Higher social welfare means that more AIGC users and AIGC service providers are participating in the AIGC network and providing high-value AIGC services within the network.
- Revenue: Providers of AIGC use a large amount of computing and energy resources to provide AIGC, which

may be offset by revenue from AIGC users. The higher the revenue, the more the AIGC service provider can be motivated to improve the AIGC service to a higher quality.
- Economic properties: In AIGC networks, AIGC providers and users should be risk-neutral, which indicates the incentive mechanisms should satisfy economic properties, e.g., individually rational, incentive compatible, and budget balance [202].

While edge learning has several promising benefits, the learning time for satisfactory performance and appropriate monetary incentives for resource providers are nontrivial challenges for AIGC. In [196], [203], [204], where mobile devices are connected to the edge server, the authors design the incentive mechanism for efficient edge learning. Specifically, mobile devices collect data and train private models locally with computational resources based on the price of edge servers in each training round. Then, the updated models are uploaded to the edge server and aggregated to minimize the global loss function. Finally, the authors not only analyze the optimal pricing strategy but also use Deep Reinforcement Learning to learn the pricing strategy to obtain the optimal solution in each round in a dynamic environment and with incomplete information. In the absence of prior knowledge, the DRL agent can learn from experience to find the optimal pricing strategy that balances payment and training time. To extend [196] to long-term incentive provisioning, the authors in [197] propose a long-term incentive mechanism for edge learning frameworks. To obtain the optimal short-term and long-term pricing strategies, the hierarchical deep reinforcement learning algorithm is used in the framework to improve the model accuracy with budget constraints.

In the process of fine-tuning the AIGC edge, the incentives described above can be used to balance the time and adaptability of the fine-tuned AIGC model. In providing incentives to AIGC service providers, the quality of AIGC services also needs to be considered in the incentive mechanism. The authors in [198] propose a quality-aware FL framework to prevent inferior model updates from degrading the global model quality. Specifically, based on an AI model trained from historical learning results, the authors estimate the learning quality of mobile devices. To motivate participants to contribute high-quality services, the authors propose a reverse auction-based incentive mechanism under the recruitment budget of edge servers, taking into account the model quality. Finally, the authors propose an algorithm for integrating the model quality into the aggregation process and for filtering non-optimal model updates to further optimize the global learning model.

Traditionally, resource utilization is inefficient, and trading mechanisms are unfair in cloud-edge computing power trading [205] for AIGC services. To address this issue, the authors in [199] develop a general trading framework for computing power grids. As illustrated in Fig. 21, the authors solve the problem of the under-utilization of computing power with AI consumers in this framework. The computing-power trading problem is first formulated as a Stackelberg game and then solved with a profit-driven multi-agent reinforcement learning algorithm. Finally, a blockchain is designed for transaction

TABLE VII: Summary of scenarios, problems, benefits/challenges, and mathematical tools of incentive mechanism.

| Ref. | Scenarios | Problems | Benefits/Challenges | Mathematical Tools |
|---|---|---|---|---|
| [196] | Efficient edge learning | A weighted sum of training time and payment/Total payment and training time | Incentivize AIGC service providers with heterogeneous resources under the uncertainty of edge network bandwidth | Deep reinforcement learning |
| [197] | Efficient edge learning | Model accuracy, number of training rounds, time efficiency/The total price | Long-term incentive mechanism for AIGC services with long-term and short-term pricing strategies | Hierarchical deep reinforcement learning |
| [198] | Quality-aware federated learning | Model accuracy and loss reduction/Learning quality estimation and quality-aware incentive mechanism | Estimate the performance of AIGC services with privacy-preserving methods for distributing proper incentives | Reverse auction |
| [199] | Cloud-Edge computing power trading for ubiquitous AI services | Profits, resource utilization, security/Computing-power unit price | Trustworthy edge-cloud resource trading framework for AIGC services | Stackelberg game and multi-agent reinforcement learning |

security in the trading framework. In mobile AIGC networks with multiple AIGC service providers and multiple AIGC users, the Stackelberg game and its extension can still provide a valid framework for equilibrium analysis. In addition, multi-agent reinforcement learning also learns the equilibrium solution of the game by exploration and exploitation in the presence of incomplete information about the game.

### F. Security and Privacy

Mobile AIGC networks leverage a collaborative computing framework on the cloud side to provide AIGC services, utilizing a large amount of heterogeneous data and computing power [206]–[208]. When mobile users are kind, AIGC can greatly enhance their creativity and efficiency. However, malicious users can also utilize AIGC for destructive purposes, posing a threat to users in mobile edge networks. For example, AI-generated text can be used by malicious users to complete phishing emails, thus compromising the security and privacy of normal users [9]. To ensure secure AIGC services, providers must choose trusted AIGC solutions and train AI models in a secure manner while providing secure hints and answers to AIGC service users.

*1) Privacy-preserving AIGC Service Provisioning:* During the lifecycle of providing AIGC services, privacy information in large-scale datasets and user requests needs to be kept secure to prevent privacy breaches. In mobile AIGC networks, the generation and storage of data for AIGC model training occur at edge servers and mobile devices [209]. Unlike resourceful cloud data centers, edge and mobile layers have limited defense capacities against various attacks. Fortunately, several privacy-preserving distributed learning frameworks, such as FL [15], have been proposed to empower privacy-preserving AIGC model fine-tuning and inference at mobile AIGC networks. In preserving user privacy in AIGC networks, FL is a distributed ML approach that allows users to transmit local models instead of data during model training [210]–[212]. Specifically, as illustrated in Fig. 20, there are two major approaches to employing FL in AIGC networks
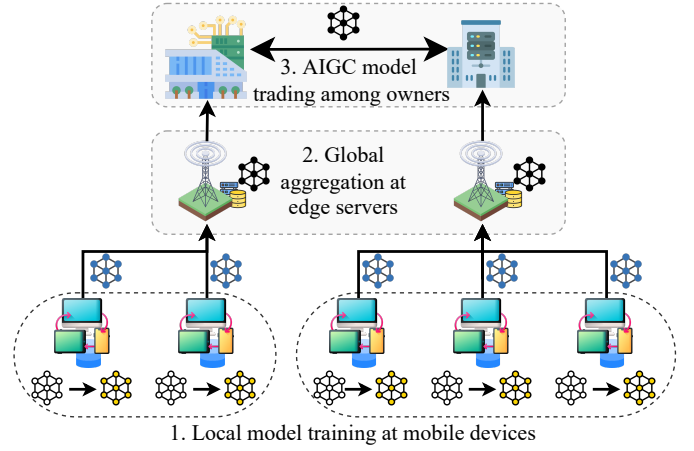


Fig. 20: Federated Learning in mobile AIGC networks, including the local model training at mobile devices, global aggregation at edge servers, and cross-server model trading.

- Secure aggregation: While FL is being learned, the mobile devices send local updates to edge servers for global aggregation. During global aggregation, authenticated encryption allows the use of secret sharing mechanisms.
- Differential privacy: Differential privacy can prevent FL servers from identifying the owners of a local update. Differential privacy is similar to secure aggregation in that it prevents FL servers from identifying owners of local updates.

Therefore, in [213], the authors propose a differential private federated generative model to synthesize representative examples of private data. With guaranteed privacy, the proposed model can solve many common data problems without human intervention. Moreover, in [214], the authors propose an FL-based generative learning scheme to improve the efficiency and robustness of GAN models. The proposed scheme is particularly effective in the presence of varying parallelism and highly skewed data distributions. To find an inherent cluster structure in users' data and unlabeled datasets, the authors

propose in [215] the unsupervised Iterative Federated Clustering algorithm, which uses generative models to deal with the statistical heterogeneity that may exist among the participants of FL. Since the centralized FL frameworks in [214], [215] might raise security concerns and risk single-point failure, the authors propose in [216] a decentralized FL framework based on a ring topology and deeply generated models. On the one hand, a method for synchronizing the ring topology can improve the communication efficiency and reliability of the system. On the other hand, generative models can solve data-related problems, such as incompleteness, low quality, insufficient quantity, and sensitivity. Finally, an InterPlanetary File System (IPFS)-based data-sharing system is developed to reduce data transmission costs and traffic congestion.

*2) Secure AIGC Service Provisioning:* Given the numerous benefits of provisioning AIGC services in mobile and edge layers, multi-tier collaboration among cloud servers, edge servers, and mobile devices enables ubiquitous AIGC service provision by heterogeneous stakeholders [217]–[220]. A trustworthy collaborative AIGC service provisioning framework must be established to provide reliable and secure AIGC services. Compared to central cloud AIGC providers, mobile and edge AIGC providers can customize AIGC services by collaborating with many user nodes while distributing data to different devices [221]. Therefore, a secure access control mechanism is required for multi-party content streaming to ensure privacy and security. However, the security of AIGC transmission cannot be ensured due to various attacks on mobile AIGC networks [222]. Fortunately, blockchain, based on distributed ledger technologies, can be utilized to explore a secure and reliable AIGC service provisioning framework and record resource and service transactions to encourage data sharing among nodes, forming a trustworthy and active mobile AIGC ecosystem [223]. As illustrated in Fig. 21, there are several benefits that blockchain brings to mobile AIGC networks [22]:

- Computing and Communication Management: Blockchain enables heterogeneous computing and communication resources to be managed securely, adaptively, and efficiently in mobile AIGC networks [224].
- Data Administration: By recording AIGC resource and service transactions in blockchain with smart contracts, data administration in mobile AIGC networks is made profitable, collaborative, and credible.
- Optimization: During optimization in AIGC services, the blockchain always provides available, complete, and secure historical data for input to optimization algorithms.

For instance, the authors in [225] propose an edge intelligence framework based on deep generative models and blockchain. To overcome the accuracy issue of the limited dataset, GAN is leveraged in the framework to synthesize training samples. Then, the output of this framework is confirmed and incentivized by smart contracts based on the proof-of-work consensus algorithm. Furthermore, the multimodal outputs of AIGC can be minted as NFTs and then recorded on the blockchain. The authors in [226] develop a conditional genera-
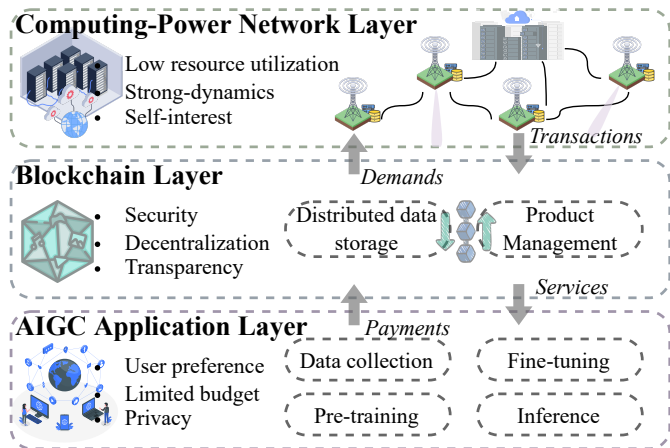


Fig. 21: Blockchain in mobile AIGC networks [199], including the AIGC application layer, blockchain layer, and computing-power network layers, for provisioning AIGC services.

tive model to synthesize new digital asset collections based on the historical transaction results of previous collections. First, the context information of NFT collections is extracted based on unsupervised learning. Based on the historical context, the newly minted collections are generated based on future token transactions. The proposed generative model can synthesize new NFT collections based on the contexts, i.e., the extracted features of previous transactions.

### G. Lessons Learned

*1) Multi-Objective Quality of AIGC Services:* In mobile AIGC networks, the quality of AIGC services is determined by several factors, including model accuracy, service latency, energy consumption, and revenue. Consequently, AIGC service providers must optimally allocate edge resources to satisfy users' multidimensional quality requirements for AIGC services [138]. Moreover, the migration of AIGC tasks and computations can enhance the reliability and efficiency of AIGC services. Notably, dynamically changing network conditions in the edge network necessitate users making online decisions to achieve load balancing and efficient use of computing resources. Attaining high-quality AIGC services requires proper considerations and practices to address the challenges discussed above, meet the quality requirements of multiple objectives, and improve user satisfaction and service quality.

*2) Edge Caching for Efficient Delivery of AIGC Services:* Edge caching plays a pivotal role in the efficient delivery of AIGC services in mobile AIGC networks. Tackling the challenges of constrained-memory edge servers, model-missing costs, and functionally equivalent models is essential for optimizing caching policies. Developing model-aware caching approaches, investigating preference-driven caching policies, and implementing principled cache designs to reduce latency and energy consumption are promising directions for enhancing the performance of mobile AIGC networks. As AI services continue to evolve, further research in caching strategies is

crucial for providing effective, personalized, and low-latency AIGC services for mobile users.

*3) Preference-aware AIGC Service Provisioning:* Offering AIGC services based on user preferences not only improves user satisfaction but also reduces service latency and resource consumption in mobile edge networks. To implement preference-based AIGC service delivery, AIGC service providers must first collect historical user data and analyze it thoroughly. In providing AIGC services, the service provider makes personalized recommendations and adjusts its strategy according to user feedback. Although user preferences play a significant role in AIGC service provision, it is essential to use and manage this information properly to protect user privacy.

*4) Life-cycle Incentive Mechanism throughout AIGC Services:* In mobile AIGC networks, the entire life cycle of AIGC services necessitates appropriate incentives for participants. A single AIGC service provider cannot provide AIGC services alone. Throughout the data collection, pre-training, fine-tuning, and inference of AIGC services, stakeholders with heterogeneous resources require reasonable incentives and must share the benefits according to their contributions. Conversely, from the users' perspective, evaluation mechanisms must be introduced. For instance, users can assess the reputation of AIGC service providers based on their transaction history to promote service optimization and improvement. Ultimately, the provisioning and transmission logs of AIGC services can also be recorded in a tamper-proof distributed ledger.

*5) Blockchain-based System Management of Mobile AIGC Networks:* Furthermore, mobile AIGC networks connect heterogeneous user devices to edge servers and cloud data centers. This uncontrolled demand for content generation introduces uncertainty and security risks into the system. Therefore, secure management and auditing methods are required to manage devices in edge environments, such as dynamically accessing, departing, and identifying IoT devices. In the traditional centralized management architecture, the risk of central node failure is unavoidable. Thus, a secure and reliable monitoring and equipment auditing system should be developed.

## VII. FUTURE RESEARCH DIRECTIONS AND OPEN ISSUES

As listed in Table VIII, in this section, we discuss future research directions and open issues from the perspectives of networking and computing, ML, and practical implementation.

### A. Networking and Computing Issues

*1) Decentralized Mobile AIGC Networks:* With the advancement of blockchain technologies [227], decentralized mobile AIGC networks can be realized based on distributed data storage, the convergence of computing and networking, and proof-of-ownership of data [223]. Such a decentralized network structure, enabled by digital identities and smart contracts, can protect AIGC users' privacy and data security. Furthermore, based on blockchain technologies, mobile AIGC networks can achieve decentralized management of the entire lifecycle of AIGC services. Therefore, future research should

investigate specific consensus mechanisms, off-chain storage frameworks, and token structures for the deployment of decentralized mobile AIGC networks.

*2) Sustainability in Mobile AIGC Networks:* In mobile AIGC networks, the pre-training, fine-tuning, and inference of generative AI models typically consume a substantial amount of computing and networking resources [26]. Hence, future research can focus on the green operations of mobile AIGC networks that provide AIGC services with minimal energy consumption and carbon emissions. To this end, effective algorithms and frameworks should be developed to operate mobile AIGC networks under dynamic service configurations, operating modes of edge nodes, and communication links. Moreover, intelligent resource management and scheduling techniques can also be proposed to balance the tradeoff between service quality and resource consumption.

High-quality data resources are also critical for the sustainability of mobile AIGC networks [228]. The performance of generative models depends not only on effective network architectures but also on the quality of training datasets [229]. However, as AIGC becomes pervasive, training datasets are gradually replaced by synthesized data that might be irrelevant to real data. Therefore, improving the quality and reliability of data in mobile AIGC networks, such as through multimodal data fusion and incremental learning technology, can further enhance the accuracy and performance of the models.

### B. Machine Learning Issues

*1) AIGC Model Compression:* As AIGC models become increasingly complex, model compression techniques are becoming more important to reduce service latency and resource consumption in provisioning AIGC services [230]. Fortunately, several techniques have been developed for AIGC model compressions, such as pruning, quantization, and knowledge distillation. First, pruning involves removing unimportant weights from the model, while quantization reduces the precision of the weights [231]. Then, knowledge distillation involves training a smaller model to mimic the larger model's behavior. Future research on AIGC model compression might continue to focus on developing and refining these techniques to improve their efficiency and effectiveness for deploying AIGC models in edge nodes and mobile devices. It is necessary to consider the limited resources of such devices and develop specialized compression techniques that can balance model size and accuracy.

*2) Privacy-preserving AIGC Services:* To provide privacy-preserving AIGC services, it is necessary to consider privacy computing techniques in both AIGC model training and inference [15]. Techniques such as differential privacy, secure multi-party computation, and homomorphic encryption can be used to protect sensitive data and prevent unauthorized access. Differential privacy involves adding noise to the data to protect individual privacy, while secure multi-party computation allows multiple parties to compute a function without revealing their inputs to one another. Homomorphic encryption enables computations to be performed on encrypted data without decryption. To successfully deploy AIGC models

TABLE VIII: A summary of future directions in mobile AIGC networks.

| Future Directions | Problems | Potential Techniques |
|---|---|---|
| Networking and Computing Issues | Decentralized Mobile AIGC Networks | Blockchain |
| | Sustainability in Mobile AIGC Networks | Green computing and communication |
| Machine Learning Issues | AIGC Model Compression | Pruning, quantization, and knowledge distillation |
| | Privacy-preserving AIGC Services | Differential privacy, secure multi-party computation, and homomorphic encryption |
| Practical Implementation Issues | Integrating AIGC and Digital Twins | Monitoring, analyzing, and predictions |
| | Immersive Streaming | AR and VR |

in edge nodes and mobile devices, the limited resources of such devices should be considered and specialized techniques that can balance privacy and performance should be developed. Additionally, concerns such as data ownership and user privacy leakage should be taken into account.

*C. Practical Implementation Issues*

*1) Integrating AIGC and Digital Twins:* Digital twins enable the maintenance of representations to monitor, analyze, and predict the status of physical entities [232]. On one hand, the integration of AIGC and digital twin technologies has the potential to significantly improve the performance of mobile AIGC networks. By creating virtual representations of physical mobile AIGC networks, service latency, and quality can be optimized through the analysis of historical data and online predictions. Furthermore, AIGC can also enhance digital twin applications by reducing the time required for designers to create simulation entities. However, several issues need to be considered during the integration of AIGC and DTs, such as efficient and secure synchronization.

*2) Immersive Streaming:* AIGC can create immersive streaming content, such as AR and VR, that can transport viewers to virtual worlds [233], which can be used in various applications such as education, entertainment, and social media. Immersive streaming can enhance the AIGC delivery process by providing a platform for viewers to interact with the generated content in real-time. However, combining AIGC and immersive streaming raises some concerns. Future research should focus on addressing the potential for biased content generation by the AIGC algorithms and the high bandwidth requirements of immersive streaming, which can cause latency issues, resulting in the degradation of the viewer's experience.

## VIII. CONCLUSIONS

In this paper, we have focused on the deployment of mobile AIGC networks, where AIGC models, services, and applications at mobile edge networks. We have discussed the background and fundamentals of generative models and the lifecycle of AIGC services at mobile AIGC networks. We have also explored AIGC-driven creative applications and use cases for mobile AIGC networks, as well as the implementation, security, and privacy challenges of deploying mobile AIGC networks. Finally, we have highlighted some future research directions and open issues for the full realization of mobile AIGC networks.

## REFERENCES

[1] E. Cetinic and J. She, "Understanding and creating art with ai: Review and outlook," *ACM Transactions on Multimedia Computing, Communications, and Applications*, vol. 18, no. 2, pp. 1–22, Feb. 2022.

[2] L.-H. Lee, Z. Lin, R. Hu, Z. Gong, A. Kumar, T. Li, S. Li, and P. Hui, "When creators meet the metaverse: A survey on computational arts," *arXiv preprint arXiv:2111.13486*, Apr. 2021, [Online]. Available: https://arxiv.org/abs/2111.13486.

[3] W. Wu, C. Zhou, M. Li, H. Wu, H. Zhou, N. Zhang, X. S. Shen, and W. Zhuang, "AI-native network slicing for 6G networks," *IEEE Wireless Communications*, vol. 29, no. 1, pp. 96–103, Apr. 2022.

[4] S. Bond-Taylor, A. Leach, Y. Long, and C. G. Willcocks, "Deep generative modelling: A comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 7327–7347, Sep. 2021.

[5] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. of the International Conference on Machine Learning*, Virtual Conference, Jul. 2021, pp. 8748–8763.

[6] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *Proc. of the International Conference on Machine Learning*, Virtual Conference, Jul. 2021, pp. 8821–8831.

[7] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, Apr. 2022, [Online]. Available: https://arxiv.org/abs/2204.06125.

[8] S. Huang, P. Grady, and GPT-3, "Generative ai: A creative new world," "Accessed Feb. 4, 2023", [Online]. Available: https://www.sequoiacap.com/article/generative-ai-a-creative-new-world/.

[9] E. Crothers, N. Japkowicz, and H. Viktor, "Machine generated text: A comprehensive survey of threat models and detection methods," *arXiv preprint arXiv:2210.07321*, Oct. 2022, [Online]. Available: https://arxiv.org/abs/2210.07321.

[10] O. AI, "Chatgpt: Optimizing language models for dialogue," "Accessed Feb. 4, 2023", [Online]. Available: https://openai.com/blog/chatgpt/.

[11] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet *et al.*, "Imagen video: High definition video generation with diffusion models," *arXiv preprint arXiv:2210.02303*, Oct. 2022, [Online]. Available: https://arxiv.org/abs/2210.02303.

[12] X. Wang, Y. Han, V. C. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of edge computing and deep learning: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 869–904, Jan. 2020.

[13] M. Westerlund, "The emergence of deepfake technology: A review," *Technology Innovation Management Review*, vol. 9, no. 11, pp. 40–53, Nov. 2019.

[14] J. Zhang and K. B. Letaief, "Mobile edge intelligence and computing for the internet of vehicles," *Proc. of the IEEE*, vol. 108, no. 2, pp. 246–261, Jun. 2019.

[15] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, Apr. 2020.

[16] M. Makhmutov, S. Varouqa, and J. A. Brow, "Survey on copyright laws about music generated by artificial intelligence," in *Proc. of the IEEE Symposium Series on Computational Intelligence*, ACT, Australia, Jan. 2020, pp. 3003–3009.

[17] M. Chen, D. Gündüz, K. Huang, W. Saad, M. Bennis, A. V. Feljan, and H. V. Poor, "Distributed learning in wireless networks: Recent progress and future challenges," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 12, pp. 3579–3605, Oct. 2021.

[18] F. Zhan, Y. Yu, R. Wu, J. Zhang, and S. Lu, "Multimodal image synthesis and editing: A survey," *arXiv preprint arXiv:2112.13592*, Dec. 2021, [Online]. Available: https://arxiv.org/abs/2112.13592.

[19] X. Shen, J. Gao, W. Wu, M. Li, C. Zhou, and W. Zhuang, "Holistic network virtualization and pervasive network intelligence for 6g," *IEEE*

*Communications Surveys & Tutorials*, vol. 24, no. 1, pp. 1–30, Dec. 2021.

[20] K. B. Letaief, Y. Shi, J. Lu, and J. Lu, "Edge artificial intelligence for 6g: Vision, enabling technologies, and applications," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 1, pp. 5–36, Nov. 2021.

[21] H. Cao, C. Tan, Z. Gao, G. Chen, P.-A. Heng, and S. Z. Li, "A survey on generative diffusion model," *arXiv preprint arXiv:2209.02646*, Sep. 2022, [Online]. Available: https://arxiv.org/abs/2209.02646.

[22] X. Wang, X. Ren, C. Qiu, Z. Xiong, H. Yao, and V. C. Leung, "Integrating edge intelligence and blockchain: What, why, and how," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 4, pp. 2193–2229, Jul. 2022.

[23] 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," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 656–700, Nov. 2023.

[24] S. Nyatsanga, T. Kucherenko, C. Ahuja, G. E. Henter, and M. Neff, "A comprehensive review of data-driven co-speech gesture generation," *arXiv preprint arXiv:2301.05339*, Jan. 2023, [Online]. Available: https://arxiv.org/abs/2301.05339.

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

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

[27] Z. Chen, J. Hu, X. Chen, J. Hu, X. Zheng, and G. Min, "Computation offloading and task scheduling for DNN-based applications in cloud-edge computing," *IEEE Access*, vol. 8, pp. 115 537–115 547, Jun. 2020.

[28] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," *ACM SIGARCH Computer Architecture News*, vol. 45, no. 1, pp. 615–629, Mar. 2017.

[29] H. Zhang and B. Di, "Intelligent omni-surfaces: Simultaneous refraction and reflection for full-dimensional wireless communications," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 4, pp. 1997–2028, Aug. 2022.

[30] D. Huang, P. Chen, R. Zeng, Q. Du, M. Tan, and C. GAN, "Location-aware graph convolutional networks for video question answering," in *Proc. of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, New York, New York, Feb. 2020, pp. 11 021–11 028.

[31] H. Zhang, B. Di, K. Bian, Z. Han, H. V. Poor, and L. Song, "Toward ubiquitous sensing and localization with reconfigurable intelligent surfaces," *Proceedings of the IEEE*, vol. 110, no. 9, pp. 1401–1422, May 2022.

[32] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. C. Leung, "Cache in the air: Exploiting content caching and delivery techniques for 5G systems," *IEEE Communications Magazine*, vol. 52, no. 2, pp. 131–139, Feb. 2014.

[33] S. Huang, H. Zhang, X. Wang, M. Chen, J. Li, and V. C. Leung, "Fine-grained spatio-temporal distribution prediction of mobile content delivery in 5G ultra-dense networks," *IEEE Transactions on Mobile Computing*, pp. 1–14, Dec. 2022.

[34] R. Gozalo-Brizuela and E. C. Garrido-Merchan, "ChatGPT is not all you need. a state of the art review of large generative ai models," *arXiv preprint arXiv:2301.04655*, Jan. 2023, [Online]. Available: https://arxiv.org/abs/2301.04655.

[35] H. Du, J. Liu, D. Niyato, J. Kang, Z. Xiong, J. Zhang, and D. I. Kim, "Attention-aware resource allocation and qoe analysis for metaverse xurllc services," *IEEE Journal on Selected Areas in Communications*, to appear, 2023.

[36] F. Tiago, F. Moreira, and T. Borges-Tiago, "Youtube videos: A destination marketing outlook," in *Proc. of the Strategic Innovative Marketing and Tourism*, Northern Aegean, Greece, May 2019, pp. 877–884.

[37] J. Krumm, N. Davies, and C. Narayanaswami, "User-generated content," *IEEE Pervasive Computing*, vol. 7, no. 4, pp. 10–11, Oct. 2008.

[38] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *arXiv preprint arXiv:2209.04747*, Sep. 2022, [Online]. Available: https://arxiv.org/abs/2209.04747.

[39] J. Oppenlaender, "Prompt engineering for text-based generative art," *arXiv preprint arXiv:2204.13988*, Apr. 2022, [Online]. Available: https://arxiv.org/abs/2204.13988.

[40] G. Marcus, E. Davis, and S. Aaronson, "A very preliminary analysis of dall-e 2," *arXiv preprint arXiv:2204.13807*, Apr. 2022, [Online]. Available: https://arxiv.org/abs/2204.13807.

[41] P.-H. Chi, P.-H. Chung, T.-H. Wu, C.-C. Hsieh, Y.-H. Chen, S.-W. Li, and H.-y. Lee, "Audio albert: A lite bert for self-supervised learning of audio representation," in *Proc. of the IEEE Spoken Language Technology Workshop*, Shenzhen, China, Jan. 2021, pp. 344–350.

[42] M. Chui, J. Manyika, M. Miremadi, N. Henke, R. Chung, P. Nel, and S. Malhotra, "Notes from the ai frontier: Insights from hundreds of use cases," *McKinsey Global Institute*, vol. 2, 2018.

[43] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Proc. of the Advances in Neural Information Processing Systems*, vol. 33, pp. 1877–1901, Dec. 2020.

[44] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *arXiv preprint arXiv:2203.02155*, Mar. 2022, [Online]. Available: https://arxiv.org/abs/2203.02155.

[45] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, Jul. 2017, [Online]. Available: https://arxiv.org/abs/1707.06347.

[46] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, "A survey for in-context learning," *arXiv preprint arXiv:2301.00234*, Jan. 2022, [Online]. Available: https://arxiv.org/abs/2301.00234.

[47] Microsoft, "Introducing the new bing," "Accessed Mar. 19, 2023", [Online]. Available: https://www.bing.com/new.

[48] J. Spataro, "Introducing microsoft 365 copilot – your copilot for work," "Accessed Mar. 19, 2023", [Online]. Available: https://blogs.microsoft.com/blog/2023/03/16/introducing-microsoft-365-copilot-your-copilot-for-work/.

[49] Y. Yang, M. Ma, H. Wu, Q. Yu, P. Zhang, X. You, J. Wu, C. Peng, T.-S. P. Yum, S. Shen *et al.*, "6g network ai architecture for everyone-centric customized services," *IEEE Network*, pp. 1–10, Jul. 2022.

[50] F. Daniel, P. Kucherbaev, C. Cappiello, B. Benatallah, and M. Allahbakhsh, "Quality control in crowdsourcing: A survey of quality attributes, assessment techniques, and assurance actions," *ACM Computing Surveys (CSUR)*, vol. 51, no. 1, pp. 1–40, Jan. 2018.

[51] H. Zhang, H. Zhang, B. Di, M. Di Renzo, Z. Han, H. V. Poor, and L. Song, "Holographic integrated sensing and communication," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 7, pp. 2114–2130, Mar. 2022.

[52] X. Deng, Y. Jiang, L. T. Yang, M. Lin, L. Yi, and M. Wang, "Data fusion based coverage optimization in heterogeneous sensor networks: A survey," *Information Fusion*, vol. 52, pp. 90–105, Dec. 2019.

[53] H. Du, J. Wang, D. Niyato, J. Kang, Z. Xiong, J. Zhang *et al.*, "Semantic communications for wireless sensing: Ris-aided encoding and self-supervised decoding," *arXiv preprint arXiv:2211.12727*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.12727.

[54] S. M. Jain, "Hugging face," in *Introduction to Transformers for NLP: With the Hugging Face Library and Models to Solve Problems*, 2022, pp. 51–67.

[55] T. Kynkäänniemi, T. Karras, S. Laine, J. Lehtinen, and T. Aila, "Improved precision and recall metric for assessing generative models," *Proc. of the Advances in Neural Information Processing Systems*, vol. 32, p. 3927–3936, Dec. 2019.

[56] D. H. Park, S. Azadi, X. Liu, T. Darrell, and A. Rohrbach, "Benchmark for compositional text-to-image synthesis," in *Proc. of the Neural Information Processing Systems Datasets and Benchmarks Track*, Virtual Conference, Dec. 2021.

[57] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, "Visual chatgpt: Talking, drawing and editing with visual foundation models," *arXiv preprint arXiv:2303.04671*, Mar. 2023, [Online]. Available: https://arxiv.org/abs/2303.04671.

[58] Y. Benny, T. Galanti, S. Benaim, and L. Wolf, "Evaluation metrics for conditional image generation," *Proc. of the International Journal of Computer Vision*, vol. 129, no. 5, pp. 1712–1731, May 2021.

[59] T. Xu, P. Zhang, Q. Huang, H. Zhang, Z. Gan, X. Huang, and X. He, "Attngan: Fine-grained text to image generation with attentional generative adversarial networks," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, Salt Lake City, Utah, Jun. 2018, pp. 1316–1324.

[60] M. F. Naeem, S. J. Oh, Y. Uh, Y. Choi, and J. Yoo, "Reliable fidelity and diversity metrics for generative models," in *Proc. of the International Conference on Machine Learning*, Virtual Conference, Nov. 2020, pp. 7176–7185.

[61] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, Oct. 2020.

[62] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," *arXiv preprint arXiv:1609.03126*, 2016.

[63] D. P. Kingma, M. Welling *et al.*, "An introduction to variational autoencoders," *Foundations and Trends® in Machine Learning*, vol. 12, no. 4, pp. 307–392, Nov. 2019.

[64] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *Proc. of the International Conference on Machine Learning*, Lille, France, Jul. 2015, pp. 1530–1538.

[65] H. T. Dinh, C. Lee, D. Niyato, and P. Wang, "A survey of mobile cloud computing: Architecture, applications, and approaches," *Wireless communications and mobile computing*, vol. 13, no. 18, pp. 1587–1611, Oct. 2013.

[66] C. Xu, Y. Ding, C. Chen, Y. Ding, W. Zhou, and S. Wen, "Personalized location privacy protection for location-based services in vehicular networks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 9, no. 10, pp. 1633–1637, Jul. 2022.

[67] M. Zhang and J. Li, "A commentary of gpt-3 in mit technology review 2021," *Fundamental Research*, vol. 1, no. 6, pp. 831–833, Feb. 2021.

[68] J. D. M.-W. C. Kenton and L. K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. of NAACL-HLT*, Minneapolis, Minnesota, Jun. 2019, pp. 4171–4186.

[69] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du *et al.*, "Lamda: Language models for dialog applications," *arXiv preprint arXiv:2201.08239*, Jan. 2022, [Online]. Available: https://arxiv.org/abs/2201.082398.

[70] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Proc. of the Advances in Neural Information Processing Systems*, p. 6000–6010, Dec. 2017.

[71] Y. Zhu, R. Kiros, R. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler, "Aligning books and movies: Towards story-like visual explanations by watching movies and reading books," in *Proc. of the IEEE International Conference on Computer Vision*, Santiago, Chile, Dec. 2015, pp. 19–27.

[72] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: A method for automatic evaluation of machine translation," in *Proc. of the 40th annual meeting of the Association for Computational Linguistics*, Philadelphia, Pennsylvania, Jul. 2002, pp. 311–318.

[73] C.-Y. Lin, "ROUGE: A package for automatic evaluation of summaries," in *Text summarization branches out*, Barcelona, Spain, Jul. 2004, pp. 74–81.

[74] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, Long Beach, CA, Jun. 2019, pp. 4401–4410.

[75] A. Brock, J. Donahue, and K. Simonyan, "Large scale GAN training for high fidelity natural image synthesis," *arXiv preprint arXiv:1809.11096*, Sep. 2018, [Online]. Available: https://arxiv.org/abs/1809.11096.

[76] A. Sauer, K. Schwarz, and A. Geiger, "Stylegan-xl: Scaling stylegan to large diverse datasets," in *Proc. of the ACM SIGGRAPH*, Virtual Conference, Jul. 2022, pp. 1–10.

[77] A. Clark, J. Donahue, and K. Simonyan, "Adversarial video generation on complex datasets," *arXiv preprint arXiv:1907.06571*, Jul. 2019, [Online]. Available: https://arxiv.org/abs/1907.06571.

[78] J. Chen, H. Guo, K. Yi, B. Li, and M. Elhoseiny, "Visualgpt: Data-efficient adaptation of pretrained language models for image captioning," in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Virtual Conference, Jun. 2022, pp. 18 030–18 040.

[79] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[80] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Proc. of the Advances in Neural Information Processing Systems*, vol. 35, pp. 36 479–36 494, Nov. 2022.

[81] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *Proc. of the International Conference on Machine Learning*, Lille, France, Jul. 2015, pp. 2256–2265.

[82] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Proc. of the Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, Dec. 2020.

[83] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.

[84] I. J. Goodfellow, "On distinguishability criteria for estimating generative models," *arXiv preprint arXiv:1412.6515*, Dec. 2014, [Online]. Available: https://arxiv.org/abs/1412.6515.

[85] A. Van Den Oord, O. Vinyals *et al.*, "Neural discrete representation learning," *Proc. of the Advances in Neural Information Processing Systems*, p. 6309–6318, Dec. 2017.

[86] L. Fei-Fei, J. Deng, and K. Li, "Imagenet: Constructing a large-scale image database," *Journal of Vision*, vol. 9, no. 8, pp. 1037–1037, Jun. 2009.

[87] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proc. of the IEEE International Conference on Computer Vision*, Santiago, Chile, Dec. 2015, pp. 3730–3738.

[88] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Proc. of the European Conference on Computer Vision*, Zurich, Switzerland, Sep. 2014, pp. 740–755.

[89] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local nash equilibrium," *Proc. of the Advances in Neural Information Processing Systems*, p. 6629–6640, Dec. 2017.

[90] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training GANs," *Proc. of the Advances in Neural Information Processing Systems*, p. 2234–2242, Dec. 2016.

[91] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, Los Alamitos, CA, Jun. 2018, pp. 586–595.

[92] A. Topirceanu, G. Barina, and M. Udrescu, "MuSeNet: Collaboration in the music artists industry," in *Proc. of the European Network Intelligence Conference*, Wroclaw, Poland, Sep. 2014, pp. 89–94.

[93] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "WaveNet: A Generative Model for Raw Audio," in *Proc. of the 9th ISCA Workshop on Speech Synthesis Workshop*, Sunnyvale, California, Sep. 2016, p. 125.

[94] Z. Borsos, R. Marinier, D. Vincent, E. Kharitonov, O. Pietquin, M. Sharifi, O. Teboul, D. Grangier, M. Tagliasacchi, and N. Zeghidour, "Audiolm: a language modeling approach to audio generation," *arXiv preprint arXiv:2209.03143*, Sep. 2022, [Online]. Available: https://arxiv.org/abs/2209.03143.

[95] C. Hawthorne, A. Stasyuk, A. Roberts, I. Simon, C.-Z. A. Huang, S. Dieleman, E. Elsen, J. Engel, and D. Eck, "Enabling factorized piano music modeling and generation with the maestro dataset," *arXiv preprint arXiv:1810.12247*, Oct. 2018, [Online]. Available: https://arxiv.org/abs/1810.12247.

[96] P. Dhariwal and A. Nichol, "Diffusion models beat GANs on image synthesis," *Proc. of the Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, Dec. 2021.

[97] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video diffusion models," *arXiv preprint arXiv:2204.03458*, Apr. 2022, [Online]. Available: https://arxiv.org/abs/2204.03458.

[98] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, Sep 2022, [Online]. Available: https://arxiv.org/abs/2209.14988.

[99] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijaya-narasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, May 2017, [Online]. Available: https://arxiv.org/abs/1705.06950.

[100] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "NeRF: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, Jan. 2021.

[101] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A lite BERT for self-supervised learning of language representations," in *Proc. of the International Conference on Learning Representations*, Addis Ababa, Ethiopia, Apr. 2019.

[102] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, "MobileBERT: a compact task-agnostic BERT for resource-limited devices," in *Proc. of the 58th Annual Meeting of the Association for Computational Linguistics*, Virtual Conference, Jul. 2020, pp. 2158–2170.

[103] A. Q. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models," in

*International Conference on Machine Learning*, Baltimore, Maryland, Jun. 2022, pp. 16 784–16 804.

[104] J. Shi, C. Wu, J. Liang, X. Liu, and N. Duan, "DiVAE: Photorealistic images synthesis with denoising diffusion decoder," *arXiv preprint arXiv:2206.00386*, 2022.

[105] M. Xu, D. Niyato, J. Kang, Z. Xiong, C. Miao, and D. I. Kim, "Wireless edge-empowered metaverse: A learning-based incentive mechanism for virtual reality," in *Proc. of IEEE International Conference on Communications (ICC)*, Seoul, South Korea, Aug. 2022, pp. 5220–5225.

[106] H. Zhang, S. Mao, D. Niyato, and Z. Han, "Location-dependent augmented reality services in wireless edge-enabled metaverse systems," *IEEE Open Journal of the Communications Society*, vol. 4, pp. 171–183, Jan. 2023.

[107] O. Gafni, A. Polyak, O. Ashual, S. Sheynin, D. Parikh, and Y. Taigman, "Make-a-scene: Scene-based text-to-image generation with human priors," in *Proc. of the 17th European Conference on Computer Vision*, Tel Aviv, Israel, 2022, pp. 89–106.

[108] A. Blattmann, R. Rombach, K. Oktay, J. Müller, and B. Ommer, "Semi-parametric neural image synthesis," in *Proc. of the Advances in Neural Information Processing Systems*, Virtual Conference, Nov. 2022.

[109] H. Du, J. Wang, D. Niyato, J. Kang, Z. Xiong, X. S. Shen, and D. I. Kim, "Exploring attention-aware network resource allocation for customized metaverse services," *IEEE Network*, pp. 1–1, 2022.

[110] W. Jin, N. Ryu, G. Kim, S.-H. Baek, and S. Cho, "Dr. 3D: Adapting 3d GANs to artistic drawings," in *Proc. of the SIGGRAPH Asia*, 2022, pp. 1–8.

[111] G. Chou, Y. Bahat, and F. Heide, "Diffusionsdf: Conditional generative modeling of signed distance functions," *arXiv preprint arXiv:2211.13757*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.13757.

[112] A. Nichol, H. Jun, P. Dhariwal, P. Mishkin, and M. Chen, "Point-e: A system for generating 3D point clouds from complex prompts," *arXiv preprint arXiv:2212.08751*, Dec. 2022, [Online]. Available: https://arxiv.org/abs/2212.08751.

[113] G. Metzer, E. Richardson, O. Patashnik, R. Giryes, and D. Cohen-Or, "Latent-nerf for shape-guided generation of 3d shapes and textures," *arXiv preprint arXiv:2211.07600*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2212.06135.

[114] X. Zeng, A. Vahdat, F. Williams, Z. Gojcic, O. Litany, S. Fidler, and K. Kreis, "LION: Latent point diffusion models for 3d shape generation," in *Advances in Neural Information Processing Systems*, Virtual Conference, Nov. 2022.

[115] M. Li, Y. Duan, J. Zhou, and J. Lu, "Diffusion-sdf: Text-to-shape via voxelized diffusion," *arXiv preprint arXiv:2212.03293*, Dec. 2022, [Online]. Available: https://arxiv.org/abs/2212.03293.

[116] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, and T.-Y. Lin, "Magic3d: High-resolution text-to-3d content creation," *arXiv preprint arXiv:2211.10440*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.10440.

[117] A. N. Wu, R. Stouffs, and F. Biljecki, "Generative adversarial networks in the built environment: A comprehensive review of the application of GANs across data types and scales," *Building and Environment*, p. 109477, Sep. 2022.

[118] H. Du, Z. Li, D. Niyato, J. Kang, Z. Xiong, D. I. Kim *et al.*, "Enabling AI-generated content (AIGC) services in wireless edge networks," *arXiv preprint arXiv:2301.03220*, Jan. 2023, [Online]. Available: https://arxiv.org/abs/2301.03220.

[119] Z. Li, M. Xu, J. Nie, J. Kang, W. Chen, and S. Xie, "NOMA-enabled cooperative computation offloading for blockchain-empowered internet of things: A learning approach," *IEEE Internet of Things Journal*, vol. 8, no. 4, pp. 2364–2378, Aug. 2020.

[120] W.-C. Fan, Y.-C. Chen, D. Chen, Y. Cheng, L. Yuan, and Y.-C. F. Wang, "Frido: Feature pyramid diffusion for complex scene image synthesis," *arXiv preprint arXiv:2208.13753*, Aug. 2022, [Online]. Available: https://arxiv.org/abs/2208.13753.

[121] H. Du, J. Wang, D. Niyato, J. Kang, Z. Xiong, and D. I. Kim, "AI-generated incentive mechanism and full-duplex semantic communications for information sharing," *arXiv preprint arXiv:2303.01896*, Mar. 2023, [Online]. Available: https://arxiv.org/abs/2303.01896.

[122] Y. Lin, H. Du, D. Niyato, J. Nie, J. Zhang, Y. Cheng, and Z. Yang, "Blockchain-aided secure semantic communication for AI-generated content in metaverse," *arXiv preprint arXiv:2301.11289*, Jan. 2023, [Online]. Available: https://arxiv.org/abs/2301.11289.

[123] M. Xu, D. Niyato, J. Chen, H. Zhang, J. Kang, Z. Xiong, S. Mao, and Z. Han, "Generative ai-empowered simulation for autonomous driving in vehicular mixed reality metaverses," *arXiv preprint*

*arXiv:2302.08418*, Feb. 2023, [Online]. Available: https://arxiv.org/abs/2302.08418.

[124] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, New Orleans, Louisiana, Jun. 2022, pp. 10 684–10 695.

[125] M. Mathias, R. Timofte, R. Benenson, and L. Van Gool, "Traffic sign recognition—how far are we from the solution?" in *Proc. of the International joint conference on Neural networks*, Dallas, Texas, Aug. 2013, pp. 1–8.

[126] M. Xu, D. Niyato, B. Wright, H. Zhang, J. Kang, Z. Xiong, S. Mao, and Z. Han, "Epvisa: Efficient auction design for real-time physical-virtual synchronization in the metaverse," *arXiv preprint arXiv:2211.06838*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.06838.

[127] M. Xu, D. Niyato, H. Zhang, J. Kang, Z. Xiong, S. Mao, and Z. Han, "Generative ai-empowered effective physical-virtual synchronization in the vehicular metaverse," *arXiv preprint arXiv:2301.07636*, Jan. 2023, [Online]. Available: https://arxiv.org/abs/2301.07636.

[128] J. Wang, H. Du, Z. Tian, D. Niyato, J. Kang, and X. Shen, "Semantic-aware sensing information transmission for metaverse: A contest theoretic approach," *IEEE Transactions on Wireless Communications*, pp. 1–1, Jan. 2023.

[129] J. Wang, H. Du, X. Yang, D. Niyato, J. Kang, and S. Mao, "Wireless sensing data collection and processing for metaverse avatar construction," *arXiv preprint arXiv:2211.12720*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.12720.

[130] W. Yang, H. Du, Z. Q. Liew, W. Y. B. Lim, Z. Xiong, D. Niyato, X. Chi, X. S. Shen, and C. Miao, "Semantic communications for future internet: Fundamentals, applications, and challenges," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 213–250, Nov. 2023.

[131] Y. Liu, H. Du, D. Niyato, J. Kang, Z. Xiong, C. Miao, A. Jamalipour *et al.*, "Blockchain-empowered lifecycle management for AI-generated content (AIGC) products in edge networks," *arXiv preprint arXiv:2303.02836*, Mar. 2023, [Online]. Available: https://arxiv.org/abs/2303.02836.

[132] C. Hu, W. Bao, D. Wang, and F. Liu, "Dynamic adaptive DNN surgery for inference acceleration on the edge," in *Proc. of the IEEE INFOCOM*, Paris, France, Apr. 2019, pp. 1423–1431.

[133] R. Zhang, K. Xiong, Y. Lu, B. Gao, P. Fan, and K. B. Letaief, "Joint coordinated beamforming and power splitting ratio optimization in mu-miso swipt-enabled hetnets: A multi-agent ddqn-based approach," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 2, pp. 677–693, Oct. 2021.

[134] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "When edge meets learning: Adaptive control for resource-constrained distributed machine learning," in *Proc. of the IEEE INFOCOM*, Honolulu, HI, Jun. 2018, pp. 63–71.

[135] K. Hsieh, A. Harlap, N. Vijaykumar, D. Konomis, G. R. Ganger, P. B. Gibbons, and O. Mutlu, "Gaia:{Geo-Distributed} machine learning approaching {LAN} speeds," in *Proc. of the 14th USENIX Symposium on Networked Systems Design and Implementation*, Boston, MA, Mar. 2017, pp. 629–647.

[136] Z. Lin, S. Bi, and Y.-J. A. Zhang, "Optimizing ai service placement and resource allocation in mobile edge intelligence systems," *IEEE Transactions on Wireless Communications*, vol. 20, no. 11, pp. 7257–7271, May 2021.

[137] X. Li, S. Bi, and H. Wang, "Optimizing resource allocation for joint ai model training and task inference in edge intelligence systems," *IEEE Wireless Communications Letters*, vol. 10, no. 3, pp. 532–536, Mar. 2020.

[138] K. Zhao, Z. Zhou, X. Chen, R. Zhou, X. Zhang, S. Yu, and D. Wu, "EdgeAdaptor: Online configuration adaption, model selection and resource provisioning for edge DNN inference serving at scale," *IEEE Transactions on Mobile Computing*, pp. 1–16, Jul. 2022.

[139] X. Tang, X. Chen, L. Zeng, S. Yu, and L. Chen, "Joint multiuser DNN partitioning and computational resource allocation for collaborative edge intelligence," *IEEE Internet of Things Journal*, vol. 8, no. 12, pp. 9511–9522, Jul. 2020.

[140] W. Y. B. Lim, J. S. Ng, Z. Xiong, J. Jin, Y. Zhang, D. Niyato, C. Leung, and C. Miao, "Decentralized edge intelligence: A dynamic resource allocation framework for hierarchical federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 3, pp. 536–550, Jul. 2021.

[141] R. Zhang, K. Xiong, Y. Lu, P. Fan, D. W. K. Ng, and K. B. Letaief, "Energy efficiency maximization in ris-assisted swipt networks with

rsma: A ppo-based approach," *IEEE Journal on Selected Areas in Communications*, pp. 1–1, Jan. 2023.

[142] G. Ditzler, M. Roveri, C. Alippi, and R. Polikar, "Learning in non-stationary environments: A survey," *IEEE Computational Intelligence Magazine*, vol. 10, no. 4, pp. 12–25, Nov. 2015.

[143] R. Zhang, K. Xiong, X. Tian, Y. Lu, P. Fan, and K. B. Letaief, "Inverse reinforcement learning meets power allocation in multi-user cellular networks," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, New York, NY, May 2022, pp. 1–2.

[144] W. Fan, Z. Chen, Y. Su, F. Wu, B. Tang, and Y. Liu, "Accuracy-based task offloading and resource allocation for edge intelligence in IoT," *IEEE Wireless Communications Letters*, vol. 11, no. 2, pp. 371–375, Nov. 2021.

[145] X. Chen, M. Li, H. Zhong, Y. Ma, and C.-H. Hsu, "DNNOff: offloading DNN-based intelligent IoT applications in mobile edge computing," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 4, pp. 2820–2829, Apr. 2021.

[146] B. Lin, Y. Huang, J. Zhang, J. Hu, X. Chen, and J. Li, "Cost-driven off-loading for DNN-based applications over cloud, edge, and end devices," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 8, pp. 5456–5466, Aug. 2019.

[147] L. Ren, Y. Laili, X. Li, and X. Wang, "Coding-based large-scale task assignment for industrial edge intelligence," *IEEE Transactions on Network Science and Engineering*, vol. 7, no. 4, pp. 2286–2297, Sep. 2019.

[148] H.-J. Jeong, I. Jeong, H.-J. Lee, and S.-M. Moon, "Computation offloading for machine learning web apps in the edge server environment," in *Proc. of the IEEE 38th International Conference on Distributed Computing Systems*, Vienna, Austria, Jul. 2018, pp. 1492–1499.

[149] X. Li, C. Sun, J. Wen, X. Wang, M. Guizani, and V. C. Leung, "Multi-user qoe enhancement: Federated multi-agent reinforcement learning for cooperative edge intelligence," *IEEE Network*, vol. 36, no. 5, pp. 144–151, Nov. 2022.

[150] Y. Zhan, S. Guo, P. Li, and J. Zhang, "A deep reinforcement learning-based offloading game in edge computing," *IEEE Transactions on Computers*, vol. 69, no. 6, pp. 883–893, Jan. 2020.

[151] W. Wu, P. Yang, W. Zhang, C. Zhou, and X. Shen, "Accuracy-guaranteed collaborative DNN inference in industrial iot via deep reinforcement learning," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 7, pp. 4988–4998, Aug. 2020.

[152] W. Zhang, D. Yang, H. Peng, W. Wu, W. Quan, H. Zhang, and X. Shen, "Deep reinforcement learning based resource management for DNN inference in industrial iot," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 8, pp. 7605–7618, Mar. 2021.

[153] R. Zhang, K. Xiong, W. Guo, X. Yang, P. Fan, and K. B. Letaief, "Q-learning-based adaptive power control in wireless RF energy harvesting heterogeneous networks," *IEEE Systems Journal*, vol. 15, no. 2, pp. 1861–1872, Sep. 2020.

[154] W. Wu, M. Li, K. Qu, C. Zhou, X. Shen, X. Zhuang, X. Li, and W. Shi, "Split learning over wireless networks: Parallel design and resource management," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 4, pp. 1051–1066, Feb. 2023.

[155] D. Saguil and A. Azim, "A layer-partitioning approach for faster execution of neural network-based embedded applications in edge networks," *IEEE Access*, vol. 8, pp. 59 456–59 469, Mar. 2020.

[156] Y. Matsubara, S. Baidya, D. Callegaro, M. Levorato, and S. Singh, "Distilled split deep neural networks for edge-assisted real-time systems," in *Proc. of the 2019 Workshop on Hot Topics in Video Analytics and Intelligent Edges*, Los Cabos, Mexico, Oct. 2019, pp. 21–26.

[157] K. Jiang, C. Sun, H. Zhou, X. Li, M. Dong, and V. C. Leung, "Intelligence-empowered mobile edge computing: Framework, issues, implementation, and outlook," *IEEE Network*, vol. 35, no. 5, pp. 74–82, Nov. 2021.

[158] C. Sun, X. Wu, X. Li, Q. Fan, J. Wen, and V. C. Leung, "Cooperative computation offloading for multi-access edge computing in 6g mobile networks via soft actor critic," *IEEE Transactions on Network Science and Engineering*, pp. 1–1, Apr. 2021.

[159] X. He, K. Wang, H. Lu, W. Xu, and S. Guo, "Edge QoE: Intelligent big data caching via deep reinforcement learning," *IEEE Network*, vol. 34, no. 4, pp. 8–13, Jul. 2020.

[160] T. Guo, R. J. Walls, and S. S. Ogden, "Edgeserve: efficient deep learning model caching at the edge," in *Proc. of the 4th ACM/IEEE Symposium on Edge Computing*, Arlington, Virginia, Nov. 2019, pp. 313–315.

[161] S. S. Ogden, G. R. Gilman, R. J. Walls, and T. Guo, "Many models at the edge: Scaling deep inference via model-level caching," in *Proc. of the IEEE International Conference on Autonomic Computing and Self-Organizing Systems*, Washington, DC, Sep. 2021, pp. 51–60.

[162] M. Xu, M. Zhu, Y. Liu, F. X. Lin, and X. Liu, "Deepcache: Principled cache for mobile deep vision," in *Proc. of the 24th Annual International Conference on Mobile Computing and Networking*, New Delhi, India, Oct. 2018, pp. 129–144.

[163] A. Fuerst and P. Sharma, "Faascache: keeping serverless computing alive with greedy-dual caching," in *Proc. of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, Virtual Conference, Mar. 2021, pp. 386–400.

[164] X.-Y. Zheng, M.-C. Lee, and Y.-W. P. Hong, "Knowledge caching for federated learning," in *2021 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2021, pp. 1–6.

[165] X. Wang, R. Li, C. Wang, X. Li, T. Taleb, and V. C. Leung, "Attention-weighted federated deep reinforcement learning for device-to-device assisted heterogeneous collaborative edge caching," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 154–169, Nov. 2020.

[166] M. Yao, L. Chen, J. Zhang, J. Huang, and J. Wu, "Loading cost-aware model caching and request routing for cooperative edge inference," in *Proc. of the IEEE International Conference on Communication*, Seoul, South Korea, May 2022, pp. 2327–2332.

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

[168] S. Xie, Y. Wu, S. Ma, M. Ding, Y. Shi, and M. Tang, "Robust information bottleneck for task-oriented communication with digital modulation," *arXiv preprint arXiv:2209.10382*, Sep. 2022, [Online]. Available: https://arxiv.org/abs/2209.10382.

[169] S. S. Ogden and T. Guo, "Mdinference: Balancing inference accuracy and latency for mobile applications," in *Proc. of the IEEE International Conference on Cloud Engineering*, NSW, Australia, Apr. 2020, pp. 28–39.

[170] M. Buckler, P. Bedoukian, S. Jayasuriya, and A. Sampson, "Eva$^2$: Exploiting temporal redundancy in live computer vision," in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, Los Angeles, California, Jun. 2018, pp. 533–546.

[171] E. Oakes, L. Yang, D. Zhou, K. Houck, T. Harter, A. Arpaci-Dusseau, and R. Arpaci-Dusseau, "{SOCK}: Rapid task provisioning with serverless-optimized containers," in *Proc. of the {USENIX} Annual Technical Conference ({USENIX}{ATC} 18)*, Boston, MA, Jul. 2018, pp. 57–70.

[172] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, "A joint learning and communications framework for federated learning over wireless networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, pp. 269–283, Oct. 2020.

[173] M. Xu, D. Niyato, Z. Yang, Z. Xiong, J. Kang, D. I. Kim, and X. Shen, "Privacy-preserving intelligent resource allocation for federated edge learning in quantum internet," *IEEE Journal of Selected Topics in Signal Processing*, vol. 17, no. 1, pp. 142–157, Nov. 2023.

[174] M. Chen, N. Shlezinger, H. V. Poor, Y. C. Eldar, and S. Cui, "Communication-efficient federated learning," *Proceedings of the National Academy of Sciences*, vol. 118, no. 17, p. e2024789118, Apr. 2021.

[175] Q. Wang, Z. Li, K. Nai, Y. Chen, and M. Wen, "Dynamic resource allocation for jointing vehicle-edge deep neural network inference," *Journal of Systems Architecture*, vol. 117, p. 102133, Aug. 2021.

[176] K. Yang, P. Sun, J. Lin, A. Boukerche, and L. Song, "A novel distributed task scheduling framework for supporting vehicular edge intelligence," in *Proc. of the IEEE 42nd International Conference on Distributed Computing Systems*, Bologna, Italy, Jul. 2022, pp. 972–982.

[177] Y. Sun, B. Xie, S. Zhou, and Z. Niu, "Meet: Mobility-enhanced edge intelligence for smart and green 6g networks," *IEEE Communications Magazine*, vol. 61, no. 1, pp. 64–70, Oct. 2023.

[178] D. Wang, B. Song, P. Lin, F. R. Yu, X. Du, and M. Guizani, "Resource management for edge intelligence (ei)-assisted iov using quantum-inspired reinforcement learning," *IEEE Internet of Things Journal*, vol. 9, no. 14, pp. 12 588–12 600, Dec. 2021.

[179] V. Balasubramanian, S. Otoum, and M. Reisslein, "Venet: hybrid stacked autoencoder learning for cooperative edge intelligence in IoV," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 16 643–16 653, May 2022.

[180] C. Dong, Y. Shen, Y. Qu, K. Wang, J. Zheng, Q. Wu, and F. Wu, "UAVs as an intelligent service: Boosting edge intelligence for air-

ground integrated networks," *IEEE Network*, vol. 35, no. 4, pp. 167–175, Aug. 2021.

[181] H. Luo, T. Chen, X. Li, S. Li, C. Zhang, G. Zhao, and X. Liu, "Keepedge: A knowledge distillation empowered edge intelligence framework for visual assisted positioning in uav delivery," *IEEE Transactions on Mobile Computing*, pp. 1–1, Mar. 2022.

[182] S. Zhou, Y. Sun, Z. Jiang, and Z. Niu, "Exploiting moving intelligence: Delay-optimized computation offloading in vehicular fog networks," *IEEE Communications Magazine*, vol. 57, no. 5, pp. 49–55, May 2019.

[183] S. S. Musa, M. Zennaro, M. Libsie, and E. Pietrosemoli, "Convergence of information-centric networks and edge intelligence for iov: Challenges and future directions," *Future Internet*, vol. 14, no. 7, p. 192, Jun. 2022.

[184] H. Du, D. Niyato, Y.-A. Xie, Y. Cheng, J. Kang, and D. I. Kim, "Performance analysis and optimization for jammer-aided multiantenna uav covert communication," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 10, pp. 2962–2979, Oct. 2022.

[185] J. Kang, H. Du, Z. Li, Z. Xiong, S. Ma, D. Niyato, and Y. Li, "Personalized saliency in task-oriented semantic communications: Image transmission and performance analysis," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 186–201, Nov. 2023.

[186] L. N. Huynh and E.-N. Huh, "UAV-enhanced edge intelligence: A survey," in *Proc. of the 6th International Conference on Computing Methodologies and Communication*, Erode, India, Mar. 2022, pp. 42–47.

[187] S. H. Alsamhi, F. A. Almalki, F. Afghah, A. Hawbani, A. V. Shvetsov, B. Lee, and H. Song, "Drones' edge intelligence over smart environments in b5g: Blockchain and federated learning synergy," *IEEE Transactions on Green Communications and Networking*, vol. 6, no. 1, pp. 295–312, Dec. 2021.

[188] Z. Wang, Y. Zhou, Y. Shi, and W. Zhuang, "Interference management for over-the-air federated learning in multi-cell wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 8, pp. 2361–2377, Jun. 2022.

[189] Z. Wang, J. Qiu, Y. Zhou, Y. Shi, L. Fu, W. Chen, and K. B. Letaief, "Federated learning via intelligent reflecting surface," *IEEE Transactions on Wireless Communications*, vol. 21, no. 2, pp. 808–822, Jul. 2021.

[190] W. Quan, N. Cheng, M. Qin, H. Zhang, H. A. Chan, and X. Shen, "Adaptive transmission control for software defined vehicular networks," *IEEE Wireless Communications Letters*, vol. 8, no. 3, pp. 653–656, Nov. 2018.

[191] S. Misra and S. Bera, "Soft-VAN: Mobility-aware task offloading in software-defined vehicular network," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 2, pp. 2071–2078, Dec. 2019.

[192] Y. Sun, W. Shi, X. Huang, S. Zhou, and Z. Niu, "Edge learning with timeliness constraints: Challenges and solutions," *IEEE Communications Magazine*, vol. 58, no. 12, pp. 27–33, Dec. 2020.

[193] J. Wang, K. Zhu, and E. Hossain, "Green internet of vehicles (IoV) in the 6G era: Toward sustainable vehicular communications and networking," *IEEE Transactions on Green Communications and Networking*, vol. 6, no. 1, pp. 391–423, Nov. 2021.

[194] M. Xu, D. T. Hoang, J. Kang, D. Niyato, Q. Yan, and D. I. Kim, "Secure and reliable transfer learning framework for 6g-enabled internet of vehicles," *IEEE Wireless Communications*, vol. 29, no. 4, pp. 132–139, May 2022.

[195] M. Li, J. Gao, L. Zhao, and X. Shen, "Deep reinforcement learning for collaborative edge computing in vehicular networks," *IEEE Transactions on Cognitive Communications and Networking*, vol. 6, no. 4, pp. 1122–1135, Jun. 2020.

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

[197] Y. Liu, L. Wu, Y. Zhan, S. Guo, and Z. Hong, "Incentive-driven long-term optimization for edge learning by hierarchical reinforcement mechanism," in *Proc. of IEEE 41st International Conference on Distributed Computing Systems*, DC, USA, Jul. 2021, pp. 35–45.

[198] Y. Deng, F. Lyu, J. Ren, Y.-C. Chen, P. Yang, Y. Zhou, and Y. Zhang, "Fair: Quality-aware federated learning with precise user incentive and model aggregation," in *Proc. of the IEEE INFOCOM*, BC, Canada, May 2021, pp. 1–10.

[199] X. Ren, C. Qiu, X. Wang, Z. Han, K. Xu, H. Yao, and S. Guo, "Ai-bazaar: A cloud-edge computing power trading framework for ubiquitous ai services," *IEEE Transactions on Cloud Computing*, pp. 1–1, Aug. 2022.

[200] X. Wang, Y. Zhao, C. Qiu, Z. Liu, J. Nie, and V. C. Leung, "Infedge: A blockchain-based incentive mechanism in hierarchical federated learn-

[201] Y. Zhan, J. Zhang, Z. Hong, L. Wu, P. Li, and S. Guo, "A survey of incentive mechanism design for federated learning," *IEEE Transactions on Emerging Topics in Computing*, vol. 10, no. 2, pp. 1035–1044, Mar. 2021.

[202] X. Chen, Y. Deng, G. Zhu, D. Wang, and Y. Fang, "From resource auction to service auction: An auction paradigm shift in wireless networks," *IEEE Wireless Communications*, vol. 29, no. 2, pp. 185–191, Apr. 2022.

[203] L. Wu, S. Guo, Y. Liu, Z. Hong, Y. Zhan, and W. Xu, "Sustainable federated learning with long-term online vcg auction mechanism," in *Proc. of the IEEE 42nd International Conference on Distributed Computing Systems*. Bologna, Italy: IEEE, Jul. 2022, pp. 895–905.

[204] Y. Zhan, P. Li, Z. Qu, D. Zeng, and S. Guo, "A learning-based incentive mechanism for federated learning," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6360–6368, Jan. 2020.

[205] J. Ren, G. Yu, Y. He, and G. Y. Li, "Collaborative cloud and edge computing for latency minimization," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 5, pp. 5031–5044, Mar. 2019.

[206] Z. Tian, L. Cui, J. Liang, and S. Yu, "A comprehensive survey on poisoning attacks and countermeasures in machine learning," *ACM Computing Surveys*, vol. 55, no. 8, pp. 1–35, Dec. 2022.

[207] Q. Liu, P. Li, W. Zhao, W. Cai, S. Yu, and V. C. Leung, "A survey on security threats and defensive techniques of machine learning: A data-driven view," *IEEE Access*, vol. 6, pp. 12 103–12 117, Feb. 2018.

[208] L. Xue, J. Ni, D. Liu, X. Lin, and X. Shen, "Blockchain-based fair and fine-grained data trading with privacy preservation," *IEEE Transactions on Computers*, pp. 1–1, Mar. 2023.

[209] J. Li, Y. Meng, L. Ma, S. Du, H. Zhu, Q. Pei, and X. Shen, "A federated learning based privacy-preserving smart healthcare system," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 3, pp. 2021–2031, Jul. 2021.

[210] J. Kang, X. Li, J. Nie, Y. Liu, M. Xu, Z. Xiong, D. Niyato, and Q. Yan, "Communication-efficient and cross-chain empowered federated learning for artificial intelligence of things," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 5, pp. 2966–2977, May 2022.

[211] W. Zhang, D. Yang, W. Wu, H. Peng, N. Zhang, H. Zhang, and X. Shen, "Optimizing federated learning in distributed industrial IoT: A multi-agent approach," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 12, pp. 3688–3703, Oct. 2021.

[212] L. Cui, Y. Qu, G. Xie, D. Zeng, R. Li, S. Shen, and S. Yu, "Security and privacy-enhanced federated learning for anomaly detection in iot infrastructures," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 5, pp. 3492–3500, Aug. 2021.

[213] S. Augenstein, H. B. McMahan, D. Ramage, S. Ramaswamy, P. Kairouz, M. Chen, R. Mathews *et al.*, "Generative models for effective ml on private, decentralized datasets," *arXiv preprint arXiv:1911.06679*, Nov. 2019, [Online]. Available: https://arxiv.org/abs/1911.06679.

[214] C. Fan and P. Liu, "Federated generative adversarial learning," in *Proc. of the Pattern Recognition and Computer Vision*, Nanjing, China, Oct. 2020, pp. 3–15.

[215] J. Chung, K. Lee, and K. Ramchandran, "Federated unsupervised clustering with generative models," in *Proc. of the AAAI International Workshop on Trustable, Verifiable and Auditable Federated Learning*, 2022.

[216] Z. Wang, Y. Hu, J. Xiao, and C. Wu, "Efficient ring-topology decentralized federated learning with deep generative models for industrial artificial intelligent," *Electronics*, vol. 11, no. 10, p. 1548, May 2022.

[217] S. Shen, Y. Ren, Y. Ju, X. Wang, W. Wang, and V. C. Leung, "Edge-matrix: A resource-redefined scheduling framework for sla-guaranteed multi-tier edge-cloud computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 3, pp. 820–834, Dec. 2023.

[218] K. Gai, J. Guo, L. Zhu, and S. Yu, "Blockchain meets cloud computing: A survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2009–2030, Apr. 2020.

[219] Y. Lin, Z. Gao, H. Du, D. Niyato, J. Kang, R. Deng, and X. S. Shen, "A unified blockchain-semantic framework for wireless edge intelligence enabled web 3.0," *IEEE Wireless Communications*, pp. 1–1, Mar. 2023.

[220] Y. Lin, Z. Gao, Y. Tu, H. Du, D. Niyato, J. Kang, and H. Yang, "A blockchain-based semantic exchange framework for web 3.0 toward participatory economy," *arXiv preprint arXiv:2211.16662*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.16662.

[221] Y. Lin, Z. Gao, W. Shi, Q. Wang, H. Li, M. Wang, Y. Yang, and L. Rui, "A novel architecture combining oracle with decentralized learning for

iiot," *IEEE Internet of Things Journal*, vol. 10, no. 5, pp. 3774–3785, Mar. 2023.

[222] C. Huang, W. Wang, D. Liu, R. Lu, and X. Shen, "Blockchain-assisted personalized car insurance with privacy preservation and fraud resistance," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 3, pp. 3777–3792, Mar. 2023.

[223] M. Xu, X. Ren, D. Niyato, J. Kang, C. Qiu, Z. Xiong, X. Wang, and V. Leung, "When quantum information technologies meet blockchain in web 3.0," *arXiv preprint arXiv:2211.15941*, Nov. 2022, [Online]. Available: https://arxiv.org/abs/2211.15941.

[224] Y. Lin, J. Kang, D. Niyato, Z. Gao, and Q. Wang, "Efficient consensus and elastic resource allocation empowered blockchain for vehicular networks," *IEEE Transactions on Vehicular Technology*, pp. 1–6, Dec. 2022.

[225] K. P. Dirgantoro, J. M. Lee, and D.-S. Kim, "Generative adversarial networks based on edge computing with blockchain architecture for security system," in *Proc. of the International Conference on Artificial Intelligence in Information and Communication*, Fukuoka, Japan, Feb. 2020, pp. 039–042.

[226] W. J.-W. Tann, A. Vuputuri, and E.-C. Chang, "Predicting non-fungible token (nft) collections: A contextual generative approach," *arXiv preprint arXiv:2210.15493*, Oct. 2022, [Online]. Available: https://arxiv.org/abs/2210.154935.

[227] Y. Li, C. Chen, N. Liu, H. Huang, Z. Zheng, and Q. Yan, "A blockchain-based decentralized federated learning framework with committee consensus," *IEEE Network*, vol. 35, no. 1, pp. 234–241, Dec. 2020.

[228] H. Du, Z. Li, D. Niyato, J. Kang, Z. Xiong, H. Huang, and S. Mao, "Generative AI-aided optimization for AI-generated content (AIGC) services in edge networks," *arXiv preprint arXiv:2303.13052*, 2023.

[229] Y. Shi, Y. Zhou, D. Wen, Y. Wu, C. Jiang, and K. B. Letaief, "Task-oriented communications for 6g: Vision, principles, and technologies," *arXiv preprint arXiv:2303.10920*, Mar. 2023, [Online]. Available: https://arxiv.org/abs/2303.10920.

[230] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "Model compression and acceleration for deep neural networks: The principles, progress, and challenges," *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 126–136, Jan. 2018.

[231] Z. Li, W. Su, M. Xu, R. Yu, D. Niyato, and S. Xie, "Compact learning model for dynamic off-chain routing in blockchain-based iot," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 12, pp. 3615–3630, Oct. 2022.

[232] A. El Saddik, "Digital twins: The convergence of multimedia technologies," *IEEE MultiMedia*, vol. 25, no. 2, pp. 87–92, Aug. 2018.

[233] A. Clemm, M. T. Vega, H. K. Ravuri, T. Wauters, and F. De Turck, "Toward truly immersive holographic-type communication: Challenges and solutions," *IEEE Communications Magazine*, vol. 58, no. 1, pp. 93–99, Jan. 2020.