

Blockchain-aided Secure Semantic Communication for AI-Generated Content in Metaverse

Yijing Lin, Hongyang Du, Dusit Niyato, *Fellow, IEEE*,
Jiangtian Nie, Jiayi Zhang, Yanyu Cheng, and Zhaohui Yang

Abstract—The construction of virtual transportation networks requires massive data to be transmitted from edge devices to Virtual Service Providers (VSP) to facilitate circulations between the physical and virtual domains in Metaverse. Leveraging semantic communication for reducing information redundancy, VSPs can receive semantic data from edge devices to provide varied services through advanced techniques, e.g., AI-Generated Content (AIGC), for users to explore digital worlds. But the use of semantic communication raises a security issue because attackers could send malicious semantic data with similar semantic information but different desired content to break Metaverse services and cause wrong output of AIGC. Therefore, in this paper, we first propose a blockchain-aided semantic communication framework for AIGC services in virtual transportation networks to facilitate interactions of the physical and virtual domains among VSPs and edge devices. We illustrate a training-based targeted semantic attack scheme to generate adversarial semantic data by various loss functions. We also design a semantic defense scheme that uses the blockchain and zero-knowledge proofs to tell the difference between the semantic similarities of adversarial and authentic semantic data and to check the authenticity of semantic data transformations. Simulation results show that the proposed defense method can reduce the semantic similarity of the adversarial semantic data and the authentic ones by up to 30% compared with the attack scheme.

Index Terms—Metaverse, Blockchain, Semantic Communication, Semantic Attacks, Semantic Defenses

I. INTRODUCTION

THE word Metaverse was coined in the science-fiction novel Snow Crash [1] to describe a virtual reality society in which people utilize digital avatars to symbolize themselves and experience the world. In recent years, Metaverse has received attention from academia and industries as a novel Internet application due to the advancement of Augmented Reality (AR), Virtual Reality (VR), Artificial Intelligence (AI), and Blockchain. Specifically, extending reality (AR/VR) and AI technologies can provide users with immersive experiences and facilitate continuous data synchronization from the physical domain to the virtual domain, which is supported by data

perceived by edge devices and services driven in Metaverse. Blockchain can help the physical and virtual domains to share information and construct economic systems in a decentralized manner. Thus, the Metaverse can be viewed as the integration of multiple technologies supported by massive data interactions between the physical and virtual domains.

One of the significant advantages of Metaverse is that people can conduct experiments in the virtual world that cannot be conducted in the real world. Because the virtual world can be created based on real-world data, e.g., images, sensing data, and text, the experimental result in Metaverse can be used to guide the real world. One example is the virtual transportation networks [2], i.e., virtual environments in which users can safely train automatic driving algorithms and test vehicles. To build such virtual environments, virtual service providers (VSPs) can leverage network edge devices to capture images from the real world and then render the virtual objects. However, this process entails three challenges as follows:

- D1) How to use the data collected from the real world, e.g., images, to achieve fast virtual world building?
- D2) How to improve the efficiency of data transmission to facilitate interactions of the physical and virtual domains?
- D3) How to ensure the security of the data received by VSP to ensure that the virtual world can be accurately synchronized with the real world?

Since virtual transportation networks require extensive interactions between the physical domain and Metaverse, edge devices can be utilized to capture images of geographical landmarks. However, converting these captured images into a consistent style for the virtual world is complicated. In several virtual service designs, VSPs need to hire digital painters to pre-process images [3], which is time-consuming and costly. The boom in AI has brought alternative solutions. AI-generated content (AIGC) technology [4] allows VSP to process quickly images collected from the real world using well trained AI models (**For D1**). However, the collected data could still burden the network. For example, the authors in [5] state that a pair of sensing devices can generate 3.072 megabytes of data per second, which challenges conventional communication systems. Fortunately, semantic communication [6] is introduced to filter out irrelevant information from edge devices to reduce information redundancy of VSPs by extracting semantic data from raw data and expressing desired meanings (**For D2**). With the help of semantic communications, massive amounts of data can be circulated in the virtual transportation networks to empower Metaverse services.

Corresponding author: Hongyang Du.

Yijing Lin is with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, 100876, Beijing, China (e-mail: yjlin@bupt.edu.cn).

Hongyang Du, Dusit Niyato, Jiangtian Nie, and Yanyu Cheng are with Nanyang Technological University, 639798, Singapore (e-mail: hongyang001@e.ntu.edu.sg; dniyato@ntu.edu.sg; jnie001@e.ntu.edu.sg; yanyu.cheng@ntu.edu.sg).

Jiayi Zhang is with School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing 100044, China (jiayizhang@bjtu.edu.cn).

Zhaohui Yang is with the College of Information Science and Electronic Engineering, Zhejiang University, China (e-mail: yang zhaohui@zju.edu.cn).

However, the introduction of semantic communication brings about higher requirements of the data security. Efficient semantic data sharing should be achieved between unknown VSPs and edge devices deployed in untrusted environments. However, the extracted semantic data can be tampered to almost the same descriptors (semantic similarities) but different desired meanings. The attacker (edge device) can modify the pixels of a sunflower to make it similar to the extracted semantic data (snowy mountain) in terms of semantic similarities but visually dissimilar [7], which affects the output of the AIGC models and is difficult for VSPs to detect the difference between the adversarial and authentic semantic data. Moreover, it is hard to detect and prevent semantic data mutations for virtual transportation networks in Metaverse. Since VSPs and edge devices are distributed, it is difficult to record, trace, and verify data transformations. To solve the aforementioned problems, the blockchain and Zero-Knowledge Proof techniques can be used (**For D3**). Thus, we present a blockchain-aided semantic communication framework with AIGC services to achieve virtual transportation networks in Metaverse. Here, blockchain-aided semantic communication can facilitate data circulation and economic activities between VSPs and edge devices in a decentralized manner [8]. Targeted semantic attacks [7] are utilized to generate adversarial semantic data to improve their semantic similarities almost up to that of the extracted semantic data by training various loss functions. Zero-Knowledge Proof (ZKP) [9] is integrated into the proposed framework to process semantic data and securely guarantee correct transformations. The contributions of this paper are summarized as follows:

- We propose a blockchain-aided semantic communication framework for AIGC in virtual transportation networks that ensures the authenticity of semantic data transmitted from edge devices to VSPs to facilitate interactions of the physical and virtual domains.
- We illustrate how a training-based targeted semantic attack scheme generates adversarial semantic data (images) without revealing the authentic semantic data. The attack semantic data has almost the same semantic similarities as the authentic ones but is visually dissimilar to them.
- We, for the first time, design a blockchain and zero-knowledge proof-based semantic defense scheme to assure the authentication of semantic data. The scheme can utilize zero-knowledge proof to record the transformations of semantic data, and use blockchain to track and verify semantic data mutations.

The remainder of the paper is described as follows. Section II reviews previous works on secure semantic communication. Section III demonstrates a blockchain-aided semantic communication framework for AIGC in Metaverse. Section IV illustrates a training-based targeted semantic attack scheme. Section V designs a blockchain and zero-knowledge proof-based semantic defense scheme. The proposed mechanisms are evaluated in Section VI. Section VII concludes the paper and elaborates the future work.

II. RELATED WORK

In this paper, we consider the integration of blockchain-aided semantic communications for Metaverse, which involves multiple emerging technologies. Therefore, we divide the related work into three parts: Blockchain-aided Semantic Communications, Semantic Attacks, and Semantic Defenses.

A. Blockchain-aided Semantic Communications

Semantic communication [6] can lighten virtual transportation network burdens by transmitting relevant semantic data to VSPs after processing original data by AI technologies in edge devices. Z. Weng *et al.* [10] utilized deep learning (DL) to identify the essential speech information with higher weights for semantic communication in dynamic channel conditions. To enable edge devices to perform DL-based semantic communication tasks, H. Xie *et al.* [11] proposed a lite-distributed semantic communication framework for edge devices to transmit low-complexity texts by optimizing training processes.

Although semantic communication can help Metaverse reduce information redundancy, it can not handle the challenge that how to construct trust among unknown edge devices and VSPs to facilitate data sharing and economic activities. Blockchain [12] is a peer-to-peer network that can construct decentralized ledgers for participants to share data. The integration of blockchain and AI-based semantic communication can empower Metaverse ecosystems to carry out rich activities between the physical and virtual domains [13]. Y. Lin *et al.* [8] proposed a unified blockchain-semantic framework to enable Web 3.0 services to implement on-chain and off-chain interactions. A proof of semantic mechanism is proposed to verify semantic data before adding it to blockchain. However, they do not mention the performance indicators of semantic data. Y. Lin *et al.* [14] proposed a blockchain-based semantic exchange framework that can mint Non-Fungible Tokens (NFT) for semantic data, utilize the game theory to facilitate exchange, and introduce ZKP to enable privacy-preserving. However, they do not consider semantic attacks that may reduce exchange efficiency.

B. Semantic Attacks and Defenses

The introduction of semantic communication brings about security issues for Metaverse. However, current research on the security of semantic communication is still in its infancy. Q. Hu *et al.* [15] analyzed semantic noise that causes semantic data to express misleading meanings. To reduce the effects caused by semantic noise, they added weight perturbation to adversarial training processes, suppressed noise-related and task-unrelated features, and designed semantic similarity-based loss functions to reduce transmitting overheads. X. Luo *et al.* [16] focused on privacy leakages when sharing background knowledge. They introduced symmetric encryption to adversarial training processes to encrypt and decrypt semantic data to ensure confidentiality. H. Du *et al.* [7] focused on the semantic Internet-of Things (SIoT) and proposed new performance indicators for SIoT to quantify security issues, including semantic secrecy outage probability

and detection failure probability. They focused on image transmission-oriented semantic communication and divided semantic attacks into targeted and untargeted semantic attacks. The targeted attack can generate adversarial semantic data (images) that can be recovered to a given target requested by receivers. The adversarial semantic data has almost the same descriptors (semantic similarities), but is visually dissimilar [17]. The untargeted semantic attack can generate adversarial semantic data that minimizes semantic similarities. They do not pursue to be recovered to any target by receivers.

In this paper, we study the targeted semantic attacks and introduce ZKP to differ in semantic similarities between the adversarial and target images to protect semantic data. ZKP has been widely used in blockchain to enable privacy-preserving and authenticity. R. Song *et al.* [9] utilized NFT and ZKP to construct a traceable data exchange scheme in blockchain, which can protect data privacy and exchange fairness. Z. Wang *et al.* [18] designed a ZKP-based off-chain data feed scheme to take in off-chain sensitive data to execute the business logic of smart contract-based applications. H. Galal *et al.* [19] leveraged ZKP and smart contracts to hide details of NFTs and swap NFTs in a fair manner. Y. Fang *et al.* [20] utilized ZKP to verify the authenticity of model prediction processes without leasing private parameters of deep learning models. However, the above methods do not consider how to use ZKP to prevent targeted semantic attacks to detect adversarial semantic data.

C. Artificial Intelligence Generated Content

AIGC refers to the use of AI algorithms, natural language processing (NLP), and computer vision (CV) methods to produce a variety of media forms, including text, images, and audio. In terms of text content generation, an epoch-making AIGC application is the ChatGPT, a conversational language model developed by OpenAI [21]. ChatGPT is a type of the Generative Pre-trained Transformer (GPT) model that is trained on a huge amount of conversational data. ChatGPT is able to generate text that sounds as if it were written by a human. This makes it helpful for a variety of purposes, including chatbots, virtual assistants, and language translation. For image generation, diffusion model [22], a new class of state-of-the-art generative models, have demonstrated exceptional performance in Image Generation tasks and have surpassed the performance of generative adversarial networks (GANs) [23] in numerous tasks. Stable Diffusion, a AIGC model that is released by Stability AI, is an open-source text-to-image generator that creates amazing artwork in a matter of seconds. It operates with unprecedented speed and quality on consumer-grade GPUs. Furthermore, AIGC techniques continues to advance in the field of audio generation. The authors in [24] propose a deep learning-based method that employs contrastive representation learning and clustering to automatically derive thematic information from music pieces in the training data. The simulation results show that the proposed model is capable of generating polyphonic pop piano music with repetition and plausible variations.

One of the primary benefits of AIGC is that it can be produced at a number and speed that would be difficult for

human beings to do alone. It also permits a great degree of consistency and precision. Therefore, the development of AIGC technologies can provide a strong boost to the evolution of the Internet. Notably, despite the potential benefits, there are also some worries about the ramifications of AIGC, including the possibility for prejudice, the semantic errors in the generated content, and the blurring of the boundary between human- and machine-generated content. Therefore, it is significant to ensure the correctness of the AIGC, avoiding attacks from malicious parties that causes resource waste to affect the quality of AIGC services.

III. BLOCKCHAIN-AIDED SEMANTIC COMMUNICATION FRAMEWORK FOR AIGC IN METAVERSE

The implementation of virtual transportation networks requires the following main steps: 1) Semantic extraction and transmission for data interactions between physical and virtual domains, 2) Semantic transformation and verification, and 3) AIGC for Metaverse, as shown in Fig. 1.

A. Semantic Extraction and Transmission

Pedestrians are in danger when auto-driving models in vehicles are not trained well enough, which motivates the development of virtual transportation networks in the Metaverse. Therefore, it is necessary to digitize physical domains using data produced in the real world to simulate environments for training vehicles and drivers. VSPs can take pictures using edge devices like smartphones, cameras, and sensors in physical domains to obtain information about the weather, traffic, and geographical landmarks that can be used to train and test the detection systems of vehicles in virtual domains. Virtual domain output can be fed back into physical domains to configure vehicles for better and safer performance. Based on the simulated environment, inexperienced drivers and vehicles can practice their reactions under unfamiliar weather or traffic conditions in Metaverse in a safe way. Therefore, VSPs need to frequently interact with edge devices supported by a tremendous amount of data to construct virtual transportation networks in Metaverse to provide services, which challenges the transmission capabilities of edge devices and VSPs.

Unfortunately, conventional communication systems cannot afford such frequent interactions between the physical and virtual domains, which will cause virtual transportation networks in Metaverse to lack enough data to simulate virtual environments. Semantic communication, a completely new paradigm that extracts semantic meaning from raw data for transmission, can be utilized to resolve this challenge. Instead of transmitting original data, edge devices can extract the semantic data and transmit to VSPs. The VSPs can then use the received semantic data to generate simulation environments for drivers and autonomous vehicle training on the virtual road. For instance, edge devices, e.g., smartphones, positioned at various locations along the same road may gather photographs of traffic conditions from their perspective. To reduce information redundancy, they can utilize semantic segmentation modules to crop key components of images as semantic data [2]. Then

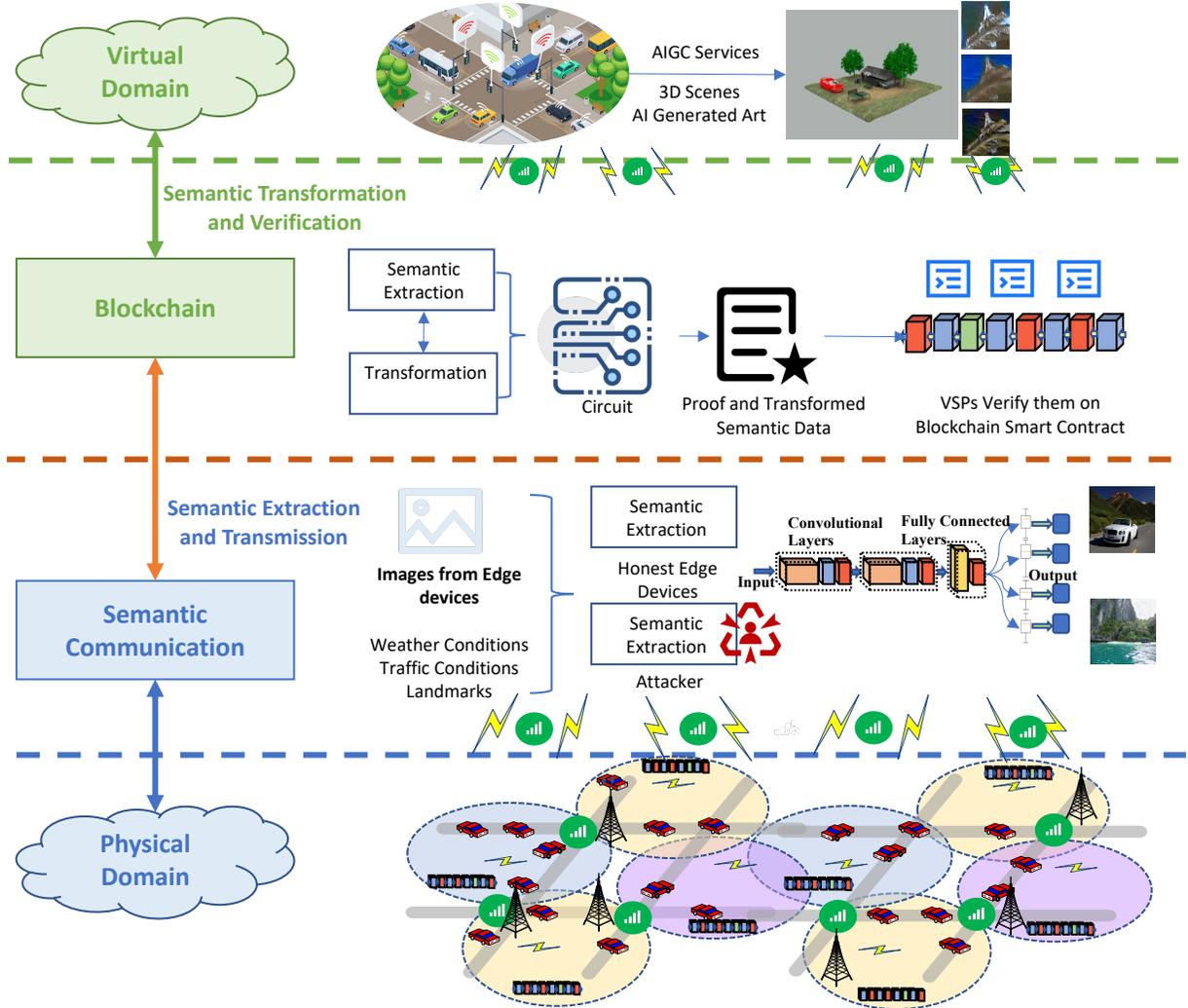


Fig. 1: Blockchain-aided Semantic Communication Framework for AIGC in Metaverse

edge devices only transmit semantic data to VSPs to report traffic conditions on roads.

However, since VSPs have to collect semantic data from multiple edge devices to train virtual transportation networks in different situations, malicious edge devices could falsify semantic data (images) and corrupt the training process, which is dangerous for pedestrians and drivers. In this paper, we study the targeted semantic attack [7, 17] in that the adversarial and authentic semantic data have almost the same semantic similarities but are totally irrelevant. The semantic similarities are calculated with the help of high-dimensional descriptors extracted by a convolutional neural network (CNN). However, malicious edge devices could train a corresponding neural network to modify some pixels in attack images to achieve similar descriptors as the authentic images to corrupt the construction of virtual transportation networks. The training-based targeted semantic attack scheme is illustrated in Section IV.

B. AIGC in Metaverse

After receiving semantic data (images) from edge devices deployed in different locations, VSPs can perceive conditions or views of landmarks, and render images by AIGC services in Metaverse. For example, semantic data of landmarks from different perspectives can be utilized by VSPs to render 3D scenes to provide users with seamless experiences. VSPs can also utilize views of landmarks to generate artworks or avatars in Metaverse. Therefore, AIGC services play an important role in virtual transportation networks to facilitate the use of data resources and the applications of Metaverse. Besides, semantic data is important for subsequent AIGC services since its quality may affect the content generated by AIGC.

However, since semantic data circulated in the Metaverse may be corrupted by the aforementioned targeted semantic attacks produced by malicious edge devices, the AIGC services may do useless work and provide users with hateful content. For example, VSPs want to collect images of famous landmarks (i.e., Eiffel Tower) while malicious edge devices may extract unrelated images of flowers to corrupt the subsequent AIGC services that renders 3D scenes with

different perspectives of the landmark. VSPs also want to perceive images of landmarks to generate artworks or avatars while attackers transmit irrelative semantic data of animals to obstacle AIGC services. Since the adversarial semantic data (images) modified by attackers have almost the same semantic similarities, it is difficult and time-consuming for VSPs to verify the authenticity of images. Therefore, it is necessary to design a mechanism to ensure the security of data transmission to protect the security of AIGC services.

C. Semantic Transformation and Verification

Malicious semantic data can affect the security of virtual transportation services in Metaverse. Modifying pixels in unrelated semantic data increases the semantic similarity score, causing the VSPs to use the wrong semantic data as input to AIGC and corrupt the virtual environment. Inspired by [17], image transformations can be performed to distinguish semantic similarities between the adversarial and authentic images. However, malicious edge devices can continue adjusting pixels in irrelative images to make them similar to transformed semantic data. Moreover, it is impossible for edge devices to transform semantic data unlimited times. A possible and practical solution is to record and verify transformations performed on semantic data.

Therefore, we propose a blockchain and zero-knowledge proof-based semantic defense scheme in Section V. The logic of transformations is recorded on the circuit produced by the zero-knowledge algorithm. Edge devices utilize extracted semantic data as inputs to generate proof of transformations and output transformed semantic data. They send the proof and semantic data to VSPs for verification. Since VSPs and edge devices are distributed in unknown Metaverse environments represented by avatars, the blockchain should be used to record and verify transformations to prevent data mutations.

IV. TRAINING-BASED TARGETED SEMANTIC ATTACK SCHEME

Although the blockchain-aided semantic communication framework for AIGC in Metaverse can reduce information redundancy and establish decentralized trust between unknown edge devices and VSPs, malicious edge devices may conduct training-based targeted semantic attacks to corrupt semantic data (images) circulated in the Metaverse services. The targeted semantic attacks refer to transmitting adversarial semantic data with almost the same semantic descriptors but visually dissimilar to the authentic one [7], which is difficult for VSPs to utilize semantic similarities (inner product of descriptors among images) evaluation to distinguish them. In this section, we illustrate the workflow of the training-based targeted semantic attack targeted semantic attacks [7][17].

Descriptor Extraction. Since extracted semantic data (images) is difficult to evaluate, edge devices can utilize a CNN to map images to high dimensional descriptors. Then VSPs can use the descriptors to distinguish semantic similarities of images. The process of descriptor extraction is illustrated as follows.

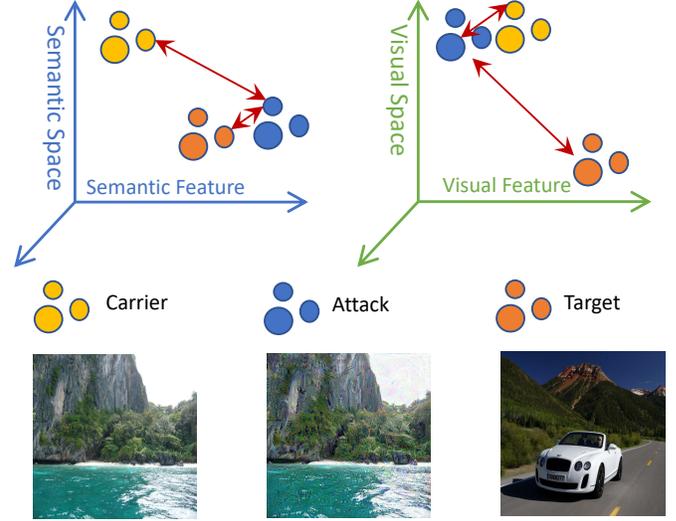


Fig. 2: Relationship among Images

Let us denote three types of semantic data produced by malicious edge devices, including the adversarial semantic data \mathbf{x}_a , the authentic one \mathbf{x}_t , and the carrier one \mathbf{x}_c . The authentic semantic data \mathbf{x}_t is extracted from original images by semantic extraction in edge devices according to requirements and associated with label $y_t = f_c(x_t)$, which has semantic similarities with \mathbf{x}_a . $f_c(\cdot)$ is utilized to classify semantic data. The carrier semantic data \mathbf{x}_c is an auxiliary image to help malicious edge devices generate \mathbf{x}_a , which is classified as $y_c = f_c(x_c) \neq y_t$ and has visual similarities with \mathbf{x}_a . The adversarial semantic data \mathbf{x}_a is produced by training networks to learn how to modify pixels, which can achieve that \mathbf{x}_a has almost the same descriptors but is visually dissimilar from the authentic one \mathbf{x}_t generated by semantic extraction, as shown in Fig. 2. Besides, \mathbf{x}_a is visually similar to \mathbf{x}_c but classified incorrectly as y_t . Therefore, the descriptor extraction process is vital for malicious edge devices to produce adversarial images \mathbf{x}_a .

The input image \mathbf{x} should be re-sampled to \mathbf{x}^s with the same dimension s to make \mathbf{x}_t and \mathbf{x}_c with the same resolution. The image \mathbf{x}^s is utilized as an input to train a Fully CNN given by $\mathbf{g}_{\mathbf{x}^s} = g(\mathbf{x}^s) : \mathbb{R}^{W \times H \times 3} \rightarrow \mathbb{R}^{w \times h \times d}$ to implement feature extraction where $W \times H \times 3$ and $w \times h \times d$ are weight, height, and channel of images. A pooling layer is utilized to map the input tensor $\mathbf{g}_{\mathbf{x}^s}$ and the descriptor $\mathbf{h}_{\mathbf{x}^s} = h(\mathbf{g}_{\mathbf{x}^s})$ by the CNN with network parameters θ , which is mapped by $h : \mathbb{R}^{w \times h \times d} \rightarrow \mathbb{R}^d$. The descriptor is the output of the pooling layer, which can be utilized to compare with other descriptors to calculate semantic similarities. The l_2 normalization is utilized to easily compare semantic similarities.

Loss Function. Considering the goal of malicious edge devices is to make \mathbf{x}_a almost the same as \mathbf{x}_t in terms of semantics but visually similar to \mathbf{x}_c , the loss function consists of the performance loss $l_{\text{TS}}(\mathbf{x}, \mathbf{x}_t)$ and the distortion loss between \mathbf{x} and \mathbf{x}_c , which can be defined as

$$L_{\text{TS}}(\mathbf{x}_c, \mathbf{x}_t; \mathbf{x}) = l_{\text{TS}}(\mathbf{x}, \mathbf{x}_t) + \lambda \|\mathbf{x} - \mathbf{x}_c\|^2, \quad (1)$$

where λ is a hyper-parameter to show the impact of the distortion loss.

Performance Loss. Let us assume that malicious edge devices can access to the network structure of the descriptor extraction [17] since it is necessary for edge devices to know the evaluation standard of VSPs. Considering different scenarios where targeted semantic attacks generate, referring to [17], we introduce the three empirical forms of the performance loss as follows.

Global Descriptor Loss is suitable that malicious edge devices even know all parameters of the network of the descriptor extraction. Thus, the performance loss function l_{TS} can be given by calculating the inner product of \mathbf{x}_a and \mathbf{x}_t as follows:

$$l_{\text{global}}(\mathbf{x}, \mathbf{x}_t) = 1 - \mathbf{h}_x^\top \mathbf{h}_{x_t}. \quad (2)$$

Activation Tensor Loss is adapted to the scenario where the outputs of the network of the descriptor extraction are the same for \mathbf{x}_a and \mathbf{x}_t before down-sampling to resolution s , which can be denoted by the mean squared difference of \mathbf{g}_x and \mathbf{g}_{x_t} as follows:

$$l_{\text{tensor}}(\mathbf{x}, \mathbf{x}_t) = \frac{\|\mathbf{g}_x - \mathbf{g}_{x_t}\|^2}{w \cdot h \cdot d}. \quad (3)$$

Activation Histogram Loss is utilized to preserve first-order statistics of activations $\|u(\mathbf{g}_x, \mathbf{b})_i\|$ per channel i to achieve identical extracted descriptors regardless of spatial information in images, which can be denoted as follows:

$$l_{\text{hist}}(\mathbf{x}, \mathbf{x}_t) = \frac{1}{d} \sum_{i=1}^d \|u(\mathbf{g}_x, \mathbf{b})_i - u(\mathbf{g}_{x_t}, \mathbf{b})_i\|, \quad (4)$$

where \mathbf{b} is the histogram bin centers.

Optimization. Our goal is to find the optimal loss function that minimizes the semantic similarities of \mathbf{x}_t and \mathbf{x} , which equivalently optimizes network parameters θ . We can use the Adam algorithm to update θ as

$$\theta_{t+1} = \theta_t - \eta \frac{\rho_t}{\sqrt{v_t + \epsilon}}, \quad (5)$$

where η is the learning rate, ρ_t and v_t are the first-order and second-order momenta of gradients, and ϵ is utilized to prevent the $\sqrt{v_t + \epsilon}$ from being zero. Therefore, the adversarial semantic data \mathbf{x}_a can be expressed by

$$\mathbf{x}_a = \arg \min_{\mathbf{x}} L_{\text{TS}}(\mathbf{x}_c, \mathbf{x}_t; \mathbf{x}), \quad (6)$$

where L_{TS} can be replaced by l_{global} , l_{tensor} , or l_{hist} according to different scenarios.

V. BLOCKCHAIN AND ZERO-KNOWLEDGE PROOF-BASED SEMANTIC DEFENSE SCHEME

Since the adversarial semantic data generated by malicious edge devices has almost the same descriptors (semantic similarity) but different desired meanings from the authentic ones produced by honest edge devices, inspired by [7][17], we utilize blockchain and zero-knowledge proof-based semantic

defense scheme to help VSPs to identify attack images transmitted in the Metaverse. Instead of submitting extracted semantic data directly, edge devices should transform or process semantic data by the bilinear interpolation algorithm [25], and utilize Zero-Knowledge Proof to record and verify transformations. The details of the proposed scheme are elaborated as follows, as shown in Fig. 3.

Transformation. The transformation of semantic data is a training-free defense method that uses visual invariance to distinguish adversarial and authentic semantic extraction. The reason is that attackers adjust some pixels to make descriptors of adversarial images similar to the authentic ones [7]. As a result, we attempt to increase visual invariance by blurring extracted images using the spatial transformation of the bilinear interpolation algorithm.

Let us assume that (x_1, y_1) , (x_1, y_2) , (x_2, y_1) , and (x_2, y_2) are four points in the extracted images. Then the targeted points (x, y) can be obtained by the spatial transformation, i.e., bilinear interpolation in the x and y directions. The spatial transformation can be considered a mapping function $f(\cdot, \cdot)$. Thus, the linear interpolation in the x direction can be derived as follows:

$$f(x, y_1) \approx \frac{x_2 - x}{x_2 - x_1} f(x_1, y_1) + \frac{x - x_1}{x_2 - x_1} f(x_2, y_1) \quad (7)$$

and

$$f(x, y_2) \approx \frac{x_2 - x}{x_2 - x_1} f(x_1, y_2) + \frac{x - x_1}{x_2 - x_1} f(x_2, y_2). \quad (8)$$

The targeted points are transformed by the linear interpolation in the y direction as follows:

$$\begin{aligned} f(x, y) &\approx \frac{y_2 - y}{y_2 - y_1} f(x, y_1) + \frac{y - y_1}{y_2 - y_1} f(x, y_2) \\ &\approx \frac{(x_2 - x)(y_2 - y)}{(x_2 - x_1)(y_2 - y_1)} f(x_1, y_1) + \frac{(x - x_1)(y_2 - y)}{(x_2 - x_1)(y_2 - y_1)} f(x_2, y_1) \\ &\quad + \frac{(x_2 - x)(y - y_1)}{(x_2 - x_1)(y_2 - y_1)} f(x_1, y_2) + \frac{(x - x_1)(y - y_1)}{(x_2 - x_1)(y_2 - y_1)} f(x_2, y_2). \end{aligned} \quad (9)$$

Although edge devices can perform transformations to obtain the blurred semantic data that can distinguish from the adversarial one, it is difficult for VSPs to verify whether the blurred semantic data is derived from authentic transformations. Malicious edge devices may modify some pixels to make descriptors of their semantic data close to that of honest edge devices after transformations. Therefore, to assure the authenticity of the blur transformation for semantic data, we utilize ZKP to record transformations and blockchain to verify them. The proposed mechanism is a tuple of 3 polynomial-time schemes after extracting a circuit mapping the transformation, including Key Generation, Proof, and Verification, which works as follows:

Extraction. The mechanism initiates the logic in a computation circuit C using the simple arithmetic expression mapping the transformation f to construct the public statement s and the private witness w [20]. The circuit implements the logic of bilinear interpolation via circom [26], a ZKP circuit compiler. The circuit can provide a relation between the inputs and outputs semantic data which can be used for verifiable

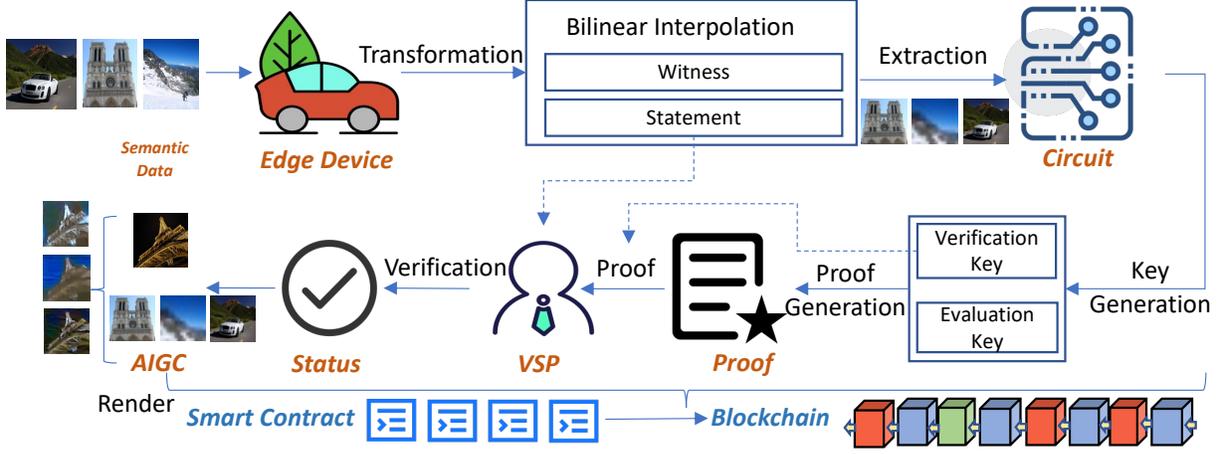


Fig. 3: Blockchain and Zero-Knowledge Proof-based Semantic Defense Scheme

computation on transformations without disclosing the inputs. The process of extraction can be illustrated as follows:

$$\text{Extract}(f) \rightarrow (s, w). \quad (10)$$

Key Generation. Edge devices take a security parameter 1^λ and the transformation circuit C as inputs to generate a common reference string crs . The common reference string crs includes an evaluation key crs.ek and a verification key crs.vk for proof and verification. The process of key generation can be expressed as follows:

$$\text{KeyGen}(1^\lambda, C) \xrightarrow{C} \text{crs}(\text{ek}, \text{vk}). \quad (11)$$

Proof. Edge devices take the evaluation key crs.ek , the statement s and the witness w related to the transformed and original semantic data, which satisfies $C(s, w) = 1$. The statement s and the witness w stand for the public and private information corresponding to the transformation relation. Thus, a zero-knowledge proof π is generated to reflect and verify the relation. The process of proof generation can be denoted as follows:

$$\text{Prove}(\text{crs.ek}, s, w) \xrightarrow{C} \pi. \quad (12)$$

Verification. VSPs utilize smart contracts deployed on blockchain to verify the authenticity of transformations and implement the business logic of blockchain-aided semantic communication for AIGC in Metaverse. Since the verification process is implemented in the blockchain, edge devices and VSPs can query verification results to construct trust in a decentralized manner. The verification key crs.vk , the statement s , and the proof π are the inputs written into smart contracts to determine whether to accept or reject the proof according to outputs. When the status of the output is 1, the verification succeeds and VSPs accept the semantic data provided by edge devices; otherwise, VSPs refuse to accept the proof corresponding to semantic data produced by malicious edge devices. The process of proof verification can be given

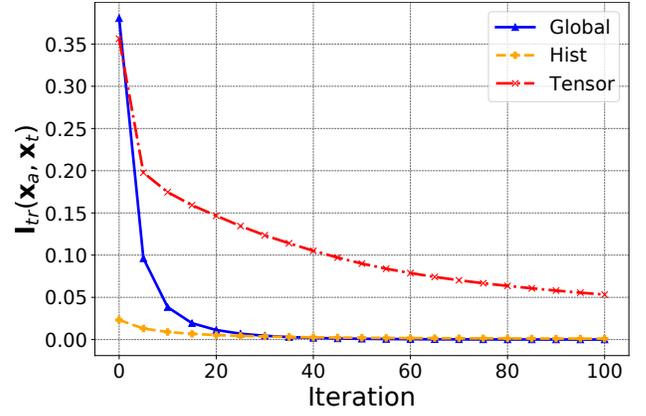


Fig. 4: Performance Loss of Attacks

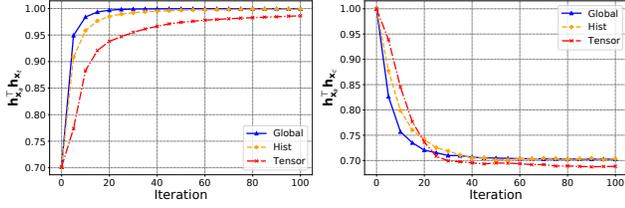
as follows:

$$\text{Verify}(\text{crs.vk}, s, \pi) \rightarrow \{0, 1\}. \quad (13)$$

The semantic defense scheme should satisfy the following properties including completeness, soundness, and zero-knowledge [9][18][19][20].

Completeness means that an honest edge device with a valid witness w can convince an honest VSP for the authenticity of transformed semantic data generated from the circuit C . If the transformed semantic data extracted by bilinear interpolation is correct and edge devices construct the proof π correctly through the proof circuit C and the evaluation key crs.ek , then the VSPs can use the verification key crs.vk generated by the proof circuit C , the proof π , and the public information (statement) s to obtain a verification result. For each $C(s, w) = 1$, the proof π generated from an honest edge device will be accepted with probability 1, which can be denoted as follows:

$$\Pr \left[\begin{array}{l} \text{KeyGen}(1^\lambda, C) \rightarrow \text{crs}(\text{ek}, \text{vk}) \\ \text{Prove}(\text{crs.ek}, s, w) \rightarrow \pi \\ \text{Verify}(\text{crs.vk}, s, \pi) = 1 \end{array} \right] = 1. \quad (14)$$



(a) Semantic similarity between the adversarial image and the authentic image (b) Semantic similarity between the adversarial image and the carrier

Fig. 5: Semantic Similarity Performance of Attacks

Soundness denotes that malicious edge devices cannot provide VSPs with authentic transformed semantic data, which means that edge devices can provide valid witnesses if they can produce valid proofs. If (s, w) is not a valid input to the proof circuit $C(s, w) = 0$, there is a polynomial time extractor Ext for a malicious edge device \mathcal{A} that makes the probabilities of $\text{Verify}(\text{crs.vk}, s, \pi^*) = 1$ are negligible. The adversarial advantage of soundness $\text{Adv}_{\mathcal{A}}^{sd}(\lambda)$ can be represented as follows:

$$\Pr \left[\begin{array}{l} \text{KeyGen}(1^\lambda, C) \rightarrow \text{crs}(\text{ek}, \text{vk}) \\ \mathcal{A}(\text{crs.ek}, s) \rightarrow \pi^* \\ \text{Ext}(\text{crs.ek}, s, \pi^*) \rightarrow w \\ \text{Verify}(\text{crs.vk}, s, \pi^*) = 1 \end{array} \right] \leq \text{negl}(\lambda). \quad (15)$$

Malicious edge devices \mathcal{A} can hardly falsify semantic data or make honest VSPs \mathcal{V} accept invalid proofs about transformations f . According to (15), VSPs cannot accept false proof π from \mathcal{A} except for a negligible probability $\text{negl}(\lambda)$.

Zero-knowledge represents that VSPs can only verify the authenticity of transformed semantic data while they cannot obtain the original semantic data unless edge devices send it to them. Let us assume that there are probabilistic polynomial time simulators S_1 , S_2 , and malicious edge devices \mathcal{A}_1 , \mathcal{A}_2 . The simulator S_1 can produce a common reference string that is used by the simulator S_2 to generate a simulated proof. Then the proposed mechanism has the zero-knowledge property if the adversarial advantage of zero-knowledge $\text{Adv}_{\mathcal{A}}^{zk}(\lambda)$ satisfies:

$$|\Pr(\text{real}) - \Pr(\text{sim})| \leq \text{negl}(\lambda), \quad (16)$$

where $\Pr(\text{real})$ and $\Pr(\text{sim})$ can be denoted as

$$\Pr(\text{real}) = \Pr \left[\begin{array}{l} \text{KeyGen}(1^\lambda, C) \rightarrow \text{crs} \\ \mathcal{A}_1(f) \rightarrow (s, w) \\ \text{Prove}(\text{crs}, s, w) \rightarrow \pi \\ \mathcal{A}_2(\text{crs}, s, w, \pi) = 1 \end{array} \right] \quad (17)$$

and

$$\Pr(\text{sim}) = \Pr \left[\begin{array}{l} S_1(1^\lambda, C) \rightarrow \text{crs} \\ \mathcal{A}_1(f) \rightarrow (s, w) \\ S_2(\text{crs}, s) \rightarrow \pi \\ \mathcal{A}_2(\text{crs}, s, w, \pi) = 1 \end{array} \right]. \quad (18)$$

VSPs who know public statements mapping with semantic data can hardly learn any knowledge about semantic data

before and after transformation. According to (16), VSPs cannot know anything about semantic data other than what can be inferred from transformation.

VI. EXPERIMENTAL EVALUATIONS

We simulate the proposed mechanism on Ubuntu 20.04LTS, Intel Xeon, 8 core, 64G memory, and 25000Mb/s with 2 Tesla V100, Go 1.19.4, Node v14.17.0, PyTorch 1.12.1, Torchvision 0.13.1, and CUDA 10.2. We perform experiments on a standard image benchmark, Revisited Paris [27], to measure the attack and defense performance among edge devices and VSPs. We exploit pre-trained networks on ImageNet [28] and AlexNet [29] to execute the attacks and defenses with 100 iterations. Unless otherwise mentioned, we use these parameters referring to [17]. The learning rate η is set to 0.01. The hyper-parameter λ that adjusts the impact of distortion loss is 0. The training process performs 100 iterations for our experiments.

A. Performance of Semantic Attack Scheme

Fig. 4 is the performance loss of the proposed semantic attack scheme with different loss functions [Global, Hist, Tensor] corresponding to the global descriptor, activation tensor, and activation histogram as the number of iterations increases. Tensor converges much slower than Global and Hist, which requires more iterations to reach its plateau.

Fig. 5 is the semantic similarity attack performance between the adversarial semantic data and the authentic or the carrier one with three loss functions [Global, Hist, Tensor]. As shown in Fig. 5 (a) and Fig. 5 (b), Global and Hist converge around the 40th iteration which is much faster than Tensor. The semantic similarity performance of attacks is consistent with Fig. 4 which shows the same trend in three loss functions.

B. Performance of Semantic Defense Scheme

Fig. 6 is the semantic similarity comparison between attack and defense schemes by displaying example images. Fig. 7 is the semantic similarity (descriptors) defense performance between the adversarial semantic data and the authentic one or the carrier one with three loss functions [Global, Hist, Tensor]. Compared with Fig. 5, the adversarial image (semantic data) can differ in semantic similarities when applying the proposed semantic defense scheme. Since the extracted semantic data is required to execute the defense scheme and the execution can be assured by zero-knowledge proof, the images can differ in descriptors. As shown in Fig. 5 (a) and Fig. 7 (a), the semantic similarity can reduce by up to 35% between the adversarial and the authentic images. Fig. 5 (b) and Fig. 7 (b) show that the semantic similarity can decrease 10% between the adversarial and the carrier images.

Fig. 8 (a) is the blockchain consensus overhead for [Attack, Authentic, Defense] types of semantic data (image) with [275 KB, 467 KB, 10 KB] data sizes when the number of nodes is [5, 50, 5]. As shown in Fig. 8 (a), the proposed scheme can

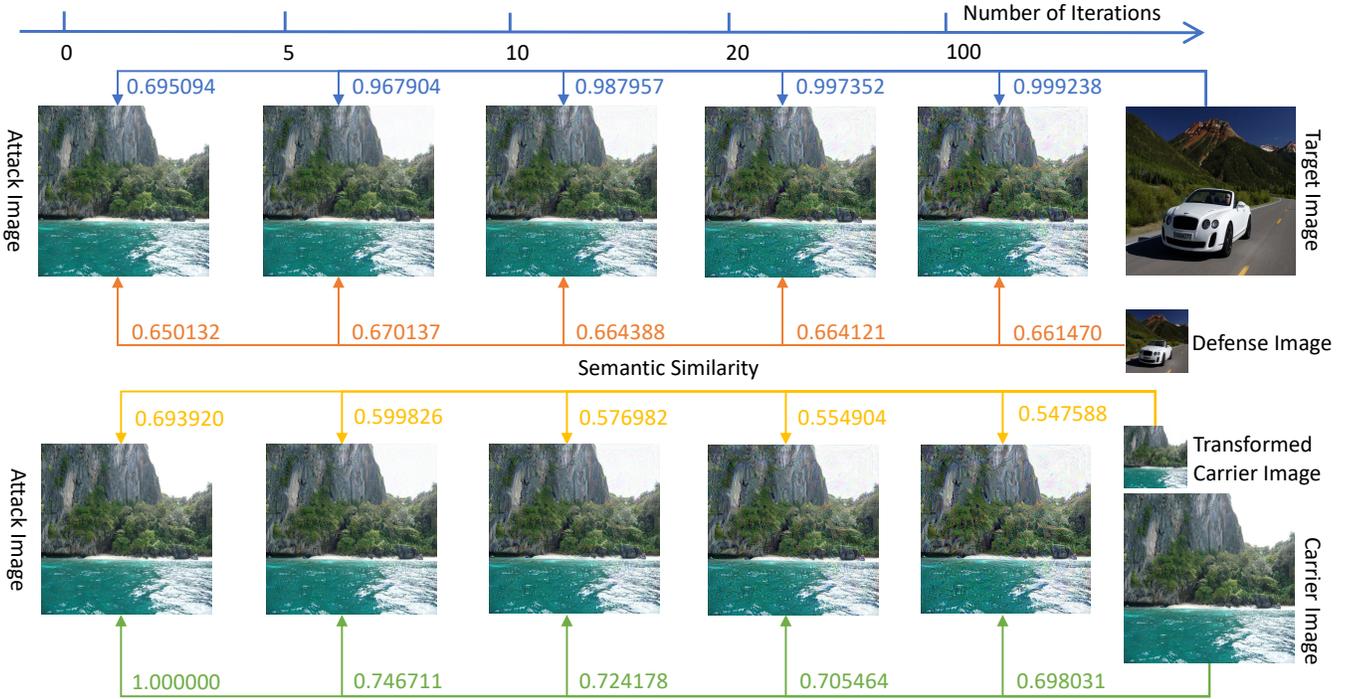
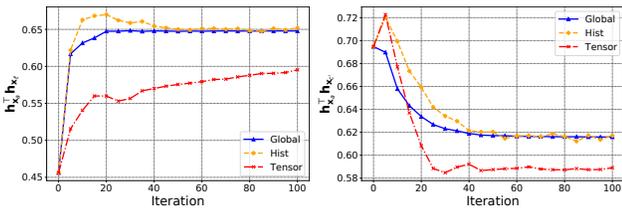


Fig. 6: Semantic Similarity Comparison Between Attack and Defense Schemes



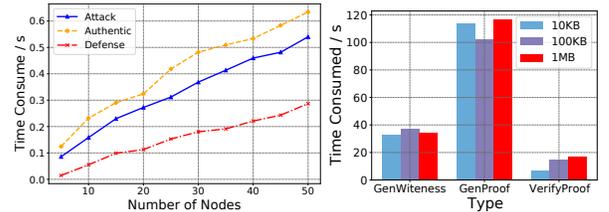
(a) Semantic similarity between the adversarial image and the authentic adversarial image and the carrier

Fig. 7: Semantic Similarity Performance of Defenses

reduce the time consumed in blockchain to improve consensus efficiency due to the attack and defense schemes cropping authentic semantic data while differing in semantic similarities according to Fig. 7. Fig. 8 (b) is the ZKP computation overhead for [GenWitness, GenProof, VerifyProof] operations with [10KB,100KB,1MB] data sizes. We can see from Fig. 8 (b) that the ZKP computation overhead depends on the GenProof operation, while the sizes of semantic data affect little on the proposed scheme. Besides, the computation overhead for verifying proofs is less than for generating proofs, which can be written into smart contracts [9][18][19].

VII. CONCLUSION AND FUTURE WORK

In this paper, we first integrated blockchain-aided semantic communication into Metaverse to support AIGC services in virtual transportation networks. We also presented a training-based targeted semantic attack scheme to illustrate potential attacks by generating adversarial semantic data with the same descriptors but different desired meanings. To prevent the



(a) Consensus

(b) ZKP Computation

Fig. 8: Consensus and Computation Overheads

above attack, we designed a blockchain and zero-knowledge proof-based semantic defense scheme that records transformations of semantic data and verifies mutations in a decentralized manner. Simulation results show that the proposed mechanism can differ in the descriptors between the adversarial semantic data and the authentic one. The defense scheme can identify malicious edge devices falsifying semantic data to corrupt AIGC services in virtual transportation networks. In future research work, we will study how to mitigate malicious edge devices and facilitate resource allocation with economic incentive mechanisms in the proposed framework.

REFERENCES

- [1] N. Stephenson, *Snow Crash: A Novel*. Spectra, 2003.
- [2] W. C. Ng, H. Du, W. Y. B. Lim, Z. Xiong, D. Niyato, and C. Miao, "Stochastic Resource Allocation for Semantic Communication-aided Virtual Transportation Networks in the Metaverse," *arXiv preprint arXiv:2208.14661*, 2022.

- [3] V. Serkova, "The digital reality: Artistic choice," in *IOP Conf. Series Materials Sci. Engin.*, vol. 940, no. 1, Jan. 2020, p. 012154.
- [4] H. Du, Z. Li, D. Niyato, J. Kang, Z. Xiong, X. Shen, and I. K. Dong, "Enabling AI-generated content (AIGC) services in wireless edge networks," *arXiv preprint arXiv:2301.03220*, 2022.
- [5] J. Wang, H. Du, Z. Tian, D. Niyato, J. Kang *et al.*, "Semantic-Aware Sensing Information Transmission for Metaverse: A Contest Theoretic Approach," *arXiv preprint arXiv:2211.12783*, 2022.
- [6] 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*, 2022.
- [7] H. Du, J. Wang, D. Niyato, J. Kang, Z. Xiong, M. Guizani, and D. I. Kim, "Rethinking Wireless Communication Security in Semantic Internet of Things," *arXiv preprint arXiv:2210.04474*, 2022.
- [8] 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," *arXiv preprint arXiv:2210.15130*, 2022.
- [9] R. Song, S. Gao, Y. Song, and B. Xiao, "ZKDET: A Traceable and Privacy-Preserving Data Exchange Scheme based on Non-Fungible Token and Zero-Knowledge," in *2022 IEEE 42nd International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2022, pp. 224–234.
- [10] Z. Weng and Z. Qin, "Semantic Communication Systems for Speech Transmission," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2434–2444, 2021.
- [11] H. Xie and Z. Qin, "A Lite Distributed Semantic Communication System for Internet of Things," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 142–153, 2020.
- [12] "Bitcoin: A Peer-to-Peer Electronic Cash System," *Decentralized Business Review*, p. 21260, 2008.
- [13] Q. Yang, Y. Zhao, H. Huang, Z. Xiong, J. Kang, and Z. Zheng, "Fusing Blockchain and AI With Metaverse: A Survey," *IEEE Open Journal of the Computer Society*, vol. 3, pp. 122–136, 2022.
- [14] 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*, 2022.
- [15] Q. Hu, G. Zhang, Z. Qin, Y. Cai, G. Yu, and G. Y. Li, "Robust Semantic Communications with Masked VQ-VAE Enabled Codebook," *arXiv preprint arXiv:2206.04011*, 2022.
- [16] X. Luo, Z. Chen, M. Tao, and F. Yang, "Encrypted Semantic Communication Using Adversarial Training for Privacy Preserving," *arXiv preprint arXiv:2209.09008*, 2022.
- [17] G. Toliás, F. Radenovic, and O. Chum, "Targeted Mismatch Adversarial Attack: Query With a Flower to Retrieve the Tower," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5037–5046.
- [18] Z. Wan, Y. Zhou, and K. Ren, "zk-AuthFeed: Protecting Data Feed to Smart Contracts with Authenticated Zero Knowledge Proof," *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [19] H. S. Galal and A. M. Youssef, "Aegis: Privacy-Preserving Market for Non-Fungible Tokens," *IEEE Transactions on Network Science and Engineering*, 2022.
- [20] Y. Fan, B. Xu, L. Zhang, J. Song, A. Zomaya, and K.-C. Li, "Validating the integrity of Convolutional Neural Network predictions based on Zero-Knowledge Proof," *Information Sciences*, 2023.
- [21] J. V. Pavlik, "Collaborating with ChatGPT: Considering the implications of generative artificial intelligence for journalism and media education," *Journal. Mass Commun. Educ.*, 2023.
- [22] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *arXiv preprint arXiv:2209.04747*, 2022.
- [23] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, "Generative adversarial networks: An overview," *IEEE Signal Processing Mag.*, vol. 35, no. 1, pp. 53–65, Jan. 2018.
- [24] Y.-J. Shih, S.-L. Wu, F. Zalkow, M. Muller, and Y.-H. Yang, "Theme transformer: Symbolic music generation with theme-conditioned transformer," *IEEE Trans. Multimedia*, to appear, 2022.
- [25] K. R. Castleman, *Digital Image Processing*. Prentice Hall Press, 1996.
- [26] Iden3, "zkSnark circuit compiler," <https://github.com/iden3/circom>, 2022.
- [27] F. Radenović, A. Iscen, G. Toliás, Y. Avrithis, and O. Chum, "Revisiting Oxford and Paris: Large-Scale Image Retrieval Benchmarking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5706–5715.
- [28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "ImageNet Large Scale Visual Recognition Challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [29] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.