# Privacy-Preserving Semantic Communication over Wiretap Channels with Learnable Differential Privacy

Weixuan Chen, *Graduate Student Member, IEEE*, Qianqian Yang, *Member, IEEE*, Shuo Shao, *Member, IEEE*,
Shunpu Tang, *Student Member, IEEE*, Zhiguo Shi, *Fellow, IEEE*, and Shui Yu, *Fellow, IEEE*

*Abstract*—While semantic communication (SemCom) improves transmission efficiency by focusing on task-relevant information, it also raises critical privacy concerns. Many existing secure SemCom approaches rely on restrictive or impractical assumptions, such as favorable channel conditions for the legitimate user or prior knowledge of the eavesdropper's model. To address these limitations, this paper proposes a novel secure SemCom framework for image transmission over wiretap channels, leveraging differential privacy (DP) to provide approximate privacy guarantees. Specifically, our approach first extracts disentangled semantic representations from source images using generative adversarial network (GAN) inversion method, and then selectively perturbs private semantic representations with approximate DP noise. Distinct from conventional DP-based protection methods, we introduce DP noise with learnable pattern, instead of traditional white Gaussian or Laplace noise, achieved through adversarial training of neural networks (NNs). This design mitigates the inherent non-invertibility of DP while effectively protecting private information. Moreover, it enables explicitly controllable security levels by adjusting the privacy budget according to specific security requirements, which is not achieved in most existing secure SemCom approaches. Experimental results demonstrate that, compared with the previous DP-based method and direct transmission, the proposed method significantly degrades the reconstruction quality for the eavesdropper, while introducing only slight degradation in task performance. Under comparable security levels, our approach achieves an LPIPS advantage of 0.06-0.29 and an FPPSR advantage of 0.10-0.86 for the legitimate user compared with the previous DP-based method.

*Index Terms*—Semantic communication, differential privacy, wiretap channel, image protection and deprotection.

## I. INTRODUCTION

### A. Backgrounds

Semantic communication (SemCom) [2] has recently emerged as a promising paradigm for future 6G networks.

Compared with conventional digital communication systems that aim to ensure bit-level accuracy, SemCom focuses on transmitting only task-relevant information. Building upon this paradigm, recent studies have demonstrated that SemCom can significantly enhance communication efficiency across various data modalities, including text [3], [4], speech [5], [6], images [7], [8], and video [9], [10]. Benefiting from these advances, SemCom is expected to enable future application scenarios such as digital twins, immersive communication, and embodied or agentic AI.

However, the SemCom approach typically eliminates traditional channel coding, which introduces redundancy that can provide a certain degree of protection for the transmitted information. Moreover, SemCom prioritizes transmitting semantically significant information while discarding as much redundancy as possible. As a result, private information, which often carries high semantic importance, becomes more vulnerable to exposure by unauthorized users. Many studies have explored secure SemCom over wiretap channels [11]–[15], which has become an important research direction and has received increasing attention in recent years.

### B. Related Works and Motivation

Researchers have explored a wide range of solutions to protect semantic information from eavesdropping. These solutions can generally be categorized into three main approaches: *adversarial training-based methods*, *encryption-based methods*, and *physical layer-based methods*. Next, we briefly review each of these approaches in the following.

*1) Adversarial training-based methods:* Marchioro *et al.* [16] proposed an adversarially trained joint source-channel coding (JSCC) framework that models the transmitter-receiver pair and the eavesdropper as players in a minimax game, aiming to confuse the eavesdropper's classifier while preserving task performance. Following a similar adversarial design, Shi *et al.* [17] introduced a controlled minimax objective that constrains the eavesdropper's performance within a specified threshold. Other studies have explored alternative JSCC variants to achieve robust adversarial learning. Erdemir *et al.* [18] proposed a VAE-based JSCC scheme trained end-to-end to capture the privacy-utility trade-off (PUT), while Zhang *et al.* [19] developed a residual-convolutional autoencoder-based JSCC framework with a SecureMSE loss, optimized to drive the eavesdropper's reconstruction toward meaningless outputs.

However, these approaches generally rely on the impractical assumption that the legitimate user has knowledge of the eavesdropper's model, and most lack mechanisms to explicitly control the level of system security.

*2) Encryption-based methods:* Some efforts have explored deep learning-based encryption. Luo *et al.* [20] proposed an encrypted SemCom system (ESCS), where the semantic message and a symmetric key are jointly processed by a neural network (NN)-based encryptor-decryptor trained adversarially to ensure confidentiality. Building upon this, Qin *et al.* [21] proposed a key generation method based on the randomness of BLEU scores in machine translation and introduced a subcarrier-level obfuscation mechanism using these semantic keys. Expanding further on cryptographic formulations, Tung *et al.* [22] proposed a framework that integrates a public-key encryption scheme based on the learning with errors (LWE) problem into a JSCC autoencoder, ensuring ciphertext indistinguishability under chosen-plaintext attacks. Moreover, Meng *et al.* [23] investigated the use of homomorphic encryption for secure SemCom. They demonstrated that semantic information can be preserved within encrypted data and designed a privacy-preserving JSCC model that supports homomorphic operations, achieving comparable task performance on both plaintext and ciphertext. However, these approaches may suffer from high computational overhead caused by sophisticated encryption schemes or from complex key management.

*3) Physical layer-based methods:* Researchers have explored introducing artificial noise into the transmitted signals in the physical layer. Chen *et al.* [24] investigated power control for artificial noise to achieve specific security levels. In their scheme, semantic information and random symbols are superposed onto a double-layered constellation. By adjusting the power allocation in superposition coding to control the symbol error probabilities (SEPs) of both users, they achieved nearly zero information leakage to the eavesdropper while maintaining reliable task performance. Building on this concept, Chen *et al.* [25] further developed a coding-enhanced jamming scheme that enables the legitimate receiver to cancel the artificial noise using shared private knowledge, while the eavesdropper cannot. Specifically, they generated private codebooks based on the shared knowledge and combined semantic information with jamming signals derived from these codebooks through superposition coding. Another approach explores the use of diffusion models. He *et al.* [26] designed and injected artificial Gaussian noise or adversarial perturbations into the transmitted signals, and employed denoising diffusion probabilistic models (DDPM) to suppress both adversarial and channel noise for the legitimate user. However, these approaches may rely on the legitimate user's channel advantage or lack theoretically provable privacy guarantees.

From the above discussion, it can be observed that existing approaches still face fundamental limitations. Most of them achieve security in specific settings or through strong assumptions, such as channel advantage, model knowledge, or pre-shared keys, which are often impractical in real systems. Moreover, current designs lack a unified mechanism to flexibly control the level of security while maintaining communication efficiency. To bridge this gap, we are motivated to design a SemCom system that provides explicitly controllable security in the challenging comparable-SNR wiretap channel scenario [27], where the eavesdropper and the legitimate receiver experience comparable channel SNRs, while eliminating the need for key exchange.

To achieve this, we are inspired to introduce differential privacy (DP) [28] into SemCom, as DP is a well-established framework that provides quantifiable privacy guarantees by injecting controlled noise to obscure sensitive information while preserving data utility. While previous studies have explored applying DP to protect source data, such as visual or semantic representations [29]–[33], directly applying DP to SemCom over wiretap channels remains challenging. This is because the inherent non-invertibility of DP, together with the presence of channel noise, often leads to severe degradation in reconstruction quality at the legitimate receiver, thereby compromising communication reliability.

## C. Contributions

In this paper, we propose a novel secure SemCom system that leverages the DP technique to enhance security when the legitimate user and the eavesdropper experience comparable channel conditions. To efficiently extract semantic information and distinguish privacy-sensitive content, we first employ generative adversarial network (GAN) inversion to obtain disentangled semantic representations. The sensitive semantic representations are selectively determined using pre-defined indices and perturbed with learnable DP noise, while the non-sensitive parts remain unchanged to preserve semantic fidelity. This pre-definition is performed once and shared between the transmitter and the legitimate receiver, enabling fine-grained protection without key exchange.

To overcome the non-invertibility problem of traditional DP mechanisms, we further propose novel NN-based DP protection and deprotection modules. The protection module generates learnable DP noise patterns through adversarial training, ensuring that the added perturbation statistically resembles genuine DP noise for privacy preservation. Meanwhile, the deprotection module at the legitimate receiver learns to recognize and mitigate these structured noise patterns, effectively restoring the protected semantics. In addition, the intensity of the injected DP noise can be adaptively adjusted according to different privacy budgets, thereby enabling an explicit trade-off between privacy protection and reconstruction quality.

The main contributions of this paper are summarized as follows:

- We investigate a novel secure SemCom framework that leverages DP to ensure communication security in a challenging comparable-SNR wiretap channel scenario. By extracting disentangled semantic representations through GAN inversion and selectively perturbing privacy-sensitive features according to pre-shared indices, the proposed system achieves fine-grained semantic protection while eliminating the need for key exchange and maintaining high reconstruction fidelity for the legitimate user.
- We propose NN-based DP protection and deprotection modules to address the non-invertibility and performance

degradation problems of conventional DP mechanisms. Through adversarial optimization, the protection module learns to generate DP noise with learnable pattern that is statistically similar to standard DP noise, while the deprotection module effectively removes it at the receiver. Moreover, by tuning the privacy budget, our method enables explicitly controllable security levels, which most existing methods lack, enabling flexible adaptation to various privacy requirements and communication scenarios.

- We conduct extensive experiments to demonstrate that, compared with the previous DP-based method [31] and SemCom without any security mechanism, the proposed framework achieves higher system security while maintaining high-fidelity reconstruction for the legitimate user and robust adaptability under varying channel conditions. In particular, under comparable security levels, our approach achieves an LPIPS improvement of 0.06-0.29 and a face privacy protection success rate (FPPSR) improvement of 0.10-0.86 for the legitimate user compared with the previous DP-based method.

## II. SYSTEM MODEL

We consider a SemCom system designed for secure image transmission over wiretap channels, as illustrated in Fig. 1. The system comprises three entities: a transmitter (Alice), a legitimate receiver (Bob), and an eavesdropper (Eve). In this system, Alice seeks to transmit a source image, denoted by $\mathbf{X}$, to Bob over a noisy wireless channel. To preserve privacy, Alice protects the sensitive portion of the image, denoted as $\mathbf{X}_{\text{private}}$, while the non-sensitive part, $\mathbf{X}_{\text{common}}$, remains unprotected. Meanwhile, an eavesdropper, Eve, passively intercepts the transmission through a channel of comparable quality, attempting to reconstruct or infer sensitive identity-related information from the intercepted signal.

Alice utilizes a transmitter to encode and protect the source image $\mathbf{X}$, extracting the semantic representation (also referred to as latent codes) to be transmitted, denoted by $\mathbf{Z}_2$:

$$\mathbf{Z}_2 = f_{\text{Alice}}\left(\mathbf{X}\right), \tag{1}$$

where $f_{\text{Alice}}$ denotes the semantic encoder. The semantic representation $\mathbf{Z}_2$ is normalized to satisfy the average power constraint $P$, and then mapped into a complex vector $\tilde{\mathbf{Z}}_2$ by pairing values into complex symbols for transmission over an AWGN channel. Bob and Eve receive noisy semantic representations, given by

$$\mathbf{Y}_1 = \tilde{\mathbf{Z}}_2 + \mathbf{n}_1, \tag{2}$$

$$\mathbf{Y}_2 = \tilde{\mathbf{Z}}_2 + \mathbf{n}_2, \tag{3}$$

where $\mathbf{n}_{1/2} \sim \mathcal{N}(0, \sigma_{1/2}^2)$ denotes Gaussian noise with zero mean and variance $\sigma_{1/2}^2$. Here, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ denote the received semantic representations at Bob and Eve, respectively. The channel SNR between Alice and the legitimate user/eavesdropper is given by

$$\text{SNR}_{\text{leg/eve}} = 10\log_{10}\left(\frac{P}{\sigma_{1/2}^2}\right) \text{ (dB)}. \tag{4}$$

In this paper, we assume that Eve experiences channel conditions comparable to those of Bob.

Both Bob and Eve attempt to reconstruct the source image $\mathbf{X}$ as accurately as possible. The received signals $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are first converted into real-valued vectors by separating the real and imaginary parts of each complex symbol. Bob decodes $\mathbf{Y}_1$ to produce the recovered image, denoted by $\hat{\mathbf{X}}_1$:

$$\hat{\mathbf{X}}_1 = f_{\text{Bob}}\left(\mathbf{Y}_1\right), \tag{5}$$

where $f_{\text{Bob}}$ denotes Bob's decoder. Similarly, Eve decodes $\mathbf{Y}_2$ into the recovered image $\hat{\mathbf{X}}_2$:

$$\hat{\mathbf{X}}_2 = f_{\text{Eve}}\left(\mathbf{Y}_2\right), \tag{6}$$

where $f_{\text{Eve}}$ denotes Eve's decoder.

This paper aims to develop a SemCom framework that fulfills two key objectives: (1) enabling Bob to accurately recover the original image, and (2) restricting the amount of sensitive semantic information that Eve can infer or reconstruct from the intercepted signals. To assess the effectiveness of the proposed system, we employ both perceptual quality and privacy protection metrics. For image reconstruction quality, we adopt the learned perceptual image patch similarity (LPIPS) metric [34], which measures the perceptual similarity between two images in the feature space of an image classification network. As a perceptual metric, LPIPS provides a reliable approximation of human visual perception. Specifically, we adopt the AlexNet-based LPIPS [35].

To assess privacy protection, we adopt the FPPSR. Specifically, we employ the ArcFace recognition system [36], which outputs a confidence score indicating whether two facial images belong to the same individual. The FPPSR is defined as the percentage of reconstructed faces identified as different from the original. Following empirical thresholds, a face is considered different if the ArcFace score falls below 0.31 [32].

## III. PROPOSED METHOD

### A. System Overview

In this section, we provide an overview of the considered SemCom system, as illustrated in Fig. 1, which includes the transmitter, and the receivers at Bob and Eve.

*1) Transmitter:* At the transmitter, Alice utilizes the inversion method of a pre-trained Semantic StyleGAN [37] $f_{\text{inv}}$ to transform the source facial image, denoted by $\mathbf{X}$, into a disentangled semantic representation $\mathbf{Z}$, i.e., $\mathbf{Z} = f_{\text{inv}}\left(\mathbf{X}\right)$. This semantic representation comprises multiple disentangled latent codes, each associated with a specific attribute (e.g., eyes, nose, or mouth).

We denote the latent codes in $\mathbf{Z}$ that require protection as the private latent codes (also referred to as private features), $\mathbf{Z}_{\text{private}}$, while the remaining ones are referred to as the common latent codes, $\mathbf{Z}_{\text{common}}$, i.e.,

$$\mathbf{Z} = \left[\mathbf{Z}_{\text{private}}, \mathbf{Z}_{\text{common}}\right]. \tag{7}$$

Then we apply the proposed privacy-preserving mechanism to protect $\mathbf{Z}_{\text{private}}$ to protect against interception by Eve. Specifically, we utilize an NN-based DP protection module to
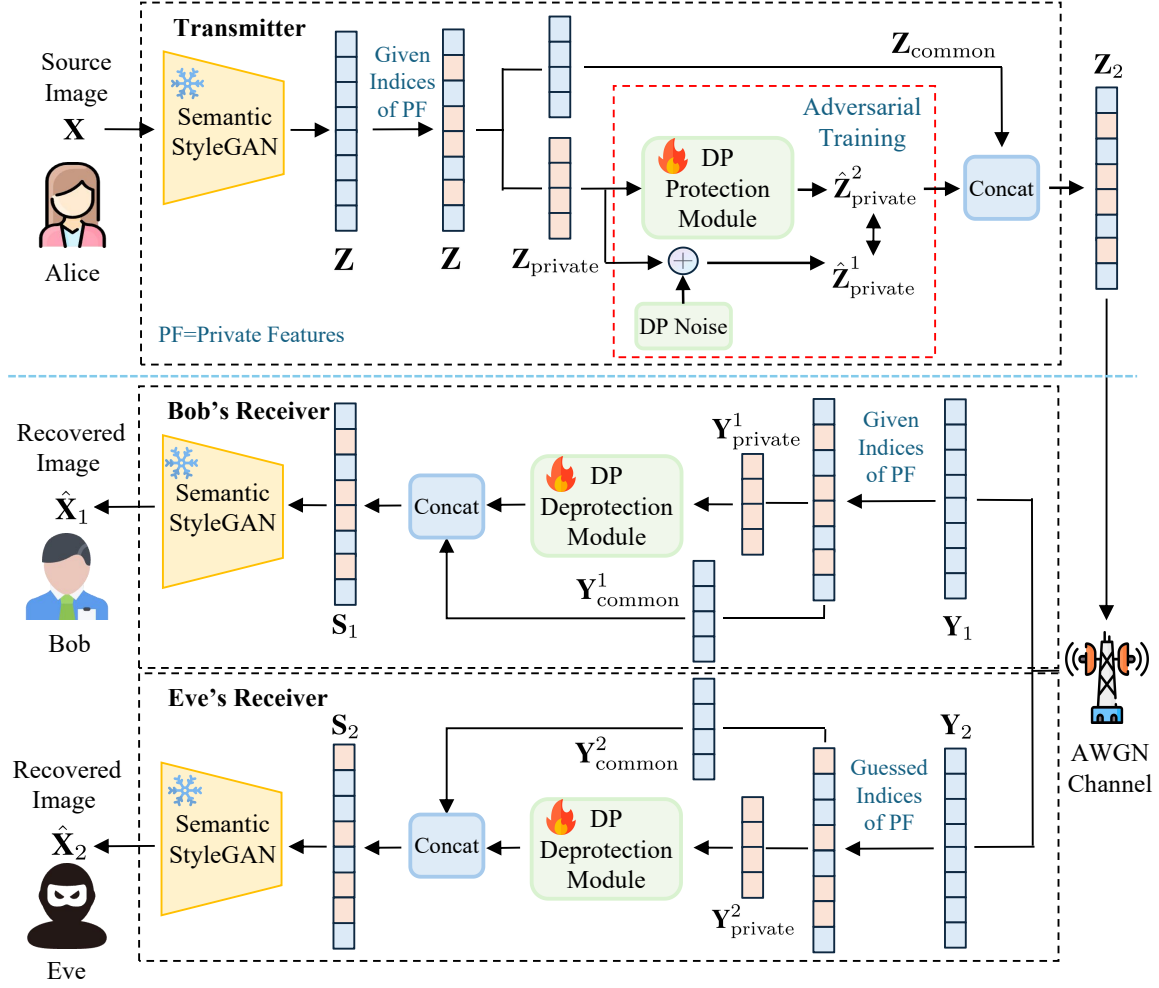
Fig. 1: The framework of our proposed secure SemCom system.

transform the private latent codes $\mathbf{Z}_{\text{private}}$ into the protected latent codes $\hat{\mathbf{Z}}^2_{\text{private}}$, i.e.,

$$\hat{\mathbf{Z}}^2_{\text{private}} = f_{\text{protection}}\left(\mathbf{Z}_{\text{private}}; \boldsymbol{\theta}^{\text{protection}}\right), \qquad (8)$$

where $f_{\text{protection}}$ denotes the NN-based DP protection module, and $\boldsymbol{\theta}^{\text{protection}}$ represents its learnable parameters. The protected private latent codes $\hat{\mathbf{Z}}^2_{\text{private}}$ are then combined with the common latent codes $\mathbf{Z}_{\text{common}}$ to form the semantic representation to be transmitted, denoted as $\mathbf{Z}_2$, i.e.,

$$\mathbf{Z}_2 = \left[\hat{\mathbf{Z}}^2_{\text{private}}, \mathbf{Z}_{\text{common}}\right]. \qquad (9)$$

*2) Bob's Receiver:* Bob knows the true indices of the private latent codes. Based on this knowledge, the real vector $\mathbf{Y}_1$ is divided into the private and common components, $\mathbf{Y}^1_{\text{private}}$ and $\mathbf{Y}^1_{\text{common}}$, respectively, following the same partitioning as at the transmitter. Next, $\mathbf{Y}^1_{\text{private}}$ is passed through an NN-based DP deprotection module to mitigate the effect of the applied protection, yielding the refined private latent codes $\hat{\mathbf{Y}}^1_{\text{private}}$:

$$\hat{\mathbf{Y}}^1_{\text{private}} = f_{\text{deprotection}}\left(\mathbf{Y}^1_{\text{private}}; \boldsymbol{\theta}^{\text{deprotection}}_1\right), \qquad (10)$$

where $f_{\text{deprotection}}$ denotes the NN-based DP deprotection module at the legitimate user, and $\boldsymbol{\theta}^{\text{deprotection}}_1$ represents its

learnable parameters. The refined private latent codes $\hat{\mathbf{Y}}^1_{\text{private}}$ are then combined with the common latent codes $\mathbf{Y}^1_{\text{common}}$ to form the complete semantic representation $\mathbf{S}_1$:

$$\mathbf{S}_1 = \left[\hat{\mathbf{Y}}^1_{\text{private}}, \mathbf{Y}^1_{\text{common}}\right]. \qquad (11)$$

Finally, Bob inputs $\mathbf{S}_1$ into a pre-trained Semantic StyleGAN generator $f_{\text{gen}}$ to reconstruct the source image, i.e.,

$$\hat{\mathbf{X}}_1 = f_{\text{gen}}\left(\mathbf{S}_1\right). \qquad (12)$$

where $\hat{\mathbf{X}}_1$ denotes the reconstructed image.

*3) Eve's Receiver:* We consider two possible scenarios for Eve. In the first scenario, Eve is unaware that the semantic representation $\tilde{\mathbf{Z}}_2$ has been protected. As a result, Eve directly inputs $\mathbf{Y}_2$ into the pre-trained Semantic StyleGAN generator to reconstruct the image $\hat{\mathbf{X}}_2$, i.e.,

$$\hat{\mathbf{X}}_2 = f_{\text{gen}}\left(\mathbf{Y}_2\right). \qquad (13)$$

In the second scenario, Eve is aware that the semantic representation $\tilde{\mathbf{Z}}_2$ has been protected and has stolen the architecture of the DP deprotection module. However, Eve does not know the trained parameters of the deprotection module, nor does she know the indices of the protected latent codes. In this case, based on the guessed indices of private latent codes, $\mathbf{Y}_2$

is divided into the private and common components, $\mathbf{Y}_{\text{private}}^2$ and $\mathbf{Y}_{\text{common}}^2$. $\mathbf{Y}_{\text{private}}^2$ is then passed through Eve's NN-based DP deprotection module to mitigate the effect of protection:

$$\hat{\mathbf{Y}}_{\text{private}}^2 = g_{\text{deprotection}}\left(\mathbf{Y}_{\text{private}}^2; \boldsymbol{\theta}_2^{\text{deprotection}}\right), \quad (14)$$

where $g_{\text{deprotection}}$ represents Eve's DP deprotection module, and $\boldsymbol{\theta}_2^{\text{deprotection}}$ denotes its learnable parameters. The refined private latent codes $\hat{\mathbf{Y}}_{\text{private}}^2$ are then combined with the common latent codes $\mathbf{Y}_{\text{common}}^2$ to form the complete semantic representation $\mathbf{S}_2$:

$$\mathbf{S}_2 = \left[\hat{\mathbf{Y}}_{\text{private}}^2, \mathbf{Y}_{\text{common}}^2\right]. \quad (15)$$

Finally, Eve inputs $\mathbf{S}_2$ into the pre-trained Semantic StyleGAN generator to reconstruct the image:

$$\hat{\mathbf{X}}_2 = f_{\text{gen}}\left(\mathbf{S}_2\right). \quad (16)$$

*4) Overall Architecture:* The proposed secure SemCom system comprises three key components: (1) a pre-trained Semantic StyleGAN generator for extracting disentangled semantic representations and reconstructing source images, (2) NN-based DP protection and deprotection modules for privacy protection and recovery, and (3) a discriminator that assists in training the DP protection module by guiding it to generate noise distributions that closely approximate genuine DP noise, while ensuring that the noise can be mitigated at the legitimate receiver. In the following, we detail each component and explain the calculation of the sensitivity $\Delta f$ in our proposed approximate DP protection mechanism.

### B. Pre-trained Semantic StyleGAN Generator

We adopt the Semantic StyleGAN generator proposed in [37] as both the encoder and decoder in our system, leveraging its bidirectional capability. Specifically, in the *forward direction*, the Semantic StyleGAN generator acts as a decoder, mapping the disentangled semantic representation to a reconstructed image at the receiver. Conversely, in the *reverse direction*, the generator functions as an encoder by performing GAN inversion at the transmitter, encoding the source image into a disentangled semantic representation. We next introduce the details of the Semantic StyleGAN generator in both the forward and reverse directions.

The framework of the Semantic StyleGAN generator is illustrated in Fig. 2. At the core of the generator are $H$ local generators, each responsible for synthesizing a specific semantic area of the face image, such as the eyes, nose, or mouth. The input semantic representation, denoted by $\mathbf{S}$, is divided into two parts: shared latent codes $\mathbf{C}^{\text{base}}$, which controls the coarse structure of the face, and $H$ local latent codes $\mathbf{C}^1 - \mathbf{C}^H$, each controlling the shape and texture of its corresponding semantic area. Each local latent code is further decomposed into a shape code and a texture code. The local generator corresponding to the $h$th semantic area takes as input both the shared latent codes and the $h$th local latent code. It is composed of modulated $1 \times 1$ convolutional layers and fully connected (FC) layers, which output a feature map $\mathbf{fm}_h$ and a pseudo-depth map $\mathbf{dm}_h$. The pseudo-depth maps are then
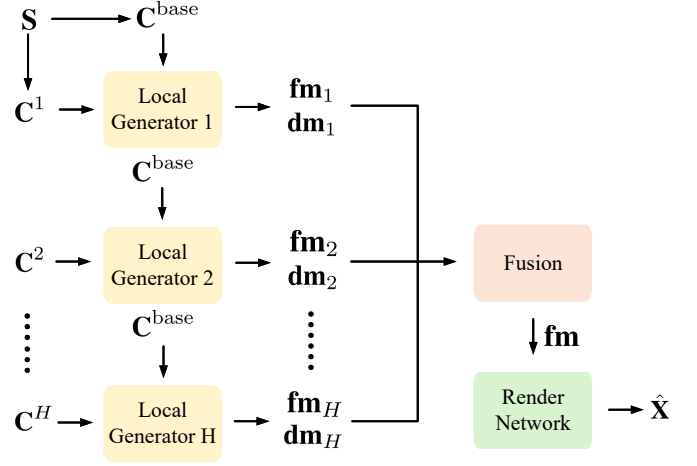


Fig. 2: The framework of the Semantic StyleGAN generator.

used to generate a coarse segmentation mask $\mathbf{m}$, which guides the fusion of all feature maps into an aggregated feature map $\mathbf{fm}$. Finally, the aggregated feature map $\mathbf{fm}$ is passed through a rendering network to produce the reconstructed image $\hat{\mathbf{X}}$.

In contrast to the forward image generation process, the reverse direction, referred to as GAN inversion, aims to encode a given source image $\mathbf{X}$ into a disentangled semantic representation $\mathbf{Z}$ that effectively captures its semantic content in the latent space. Formally, the objective of GAN inversion is to solve the following optimization problem:

$$\arg\min_{\mathbf{Z}} \text{MSE}\left(\mathbf{X}, f_{\text{gen}}(\mathbf{Z})\right), \quad (17)$$

where $f_{\text{gen}}$ represents the generator, $f_{\text{gen}}(\mathbf{Z})$ represents the image generated from $\mathbf{Z}$, and MSE represents the mean squared error (MSE). The solution can be obtained through a fixed number of gradient descent iterations.

### C. NN-based DP Protection / Deprotection Modules

*1) Discriminator and Adversarial Training:* Traditional DP protection methods are inherently non-invertible, which poses a significant challenge for reconstruction at the receiver. To overcome this limitation, we propose a novel NN-based DP protection module that generates DP noise with learnable pattern. This noise not only closely resembles genuine DP noise, thereby meeting privacy protection requirements, but also remains easily mitigated by the legitimate user through the corresponding NN-based DP deprotection module.

To protect the private latent codes using the NN-based DP protection module, the first step is to generate an intermediate variable, $\hat{\mathbf{Z}}_{\text{private}}^1$, which serves as a guiding representation during training. Specifically, $\hat{\mathbf{Z}}_{\text{private}}^1$ is obtained by adding genuine DP noise to $\mathbf{Z}_{\text{private}}$:

$$\hat{\mathbf{Z}}_{\text{private}}^1 = \mathbf{Z}_{\text{private}} + \mathbf{n}_{\text{dp}}, \quad (18)$$

where $\mathbf{n}_{\text{dp}} \sim Lap(0, \frac{\Delta f}{\epsilon})$, and the method for obtaining $\Delta f$ will be discussed in Subsection D. In the second step, the DP protection module perturbs the private latent codes under the guidance of $\hat{\mathbf{Z}}_{\text{private}}^1$, generating the perturbed latent codes
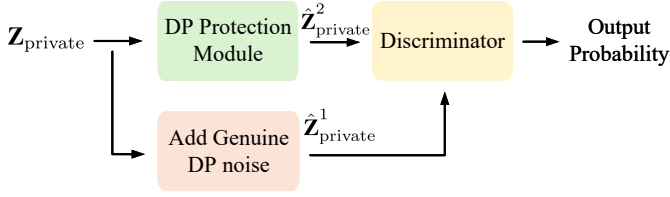
Fig. 3: Illustration of the adversarial training process for the DP protection module.



Fig. 4: The network architecture of the proposed NN-based DP protection module.

$\hat{\mathbf{Z}}^2_{\text{private}}$. It is important to emphasize that $\hat{\mathbf{Z}}^1_{\text{private}}$ serves only as an intermediate representation to guide the generation of $\hat{\mathbf{Z}}^2_{\text{private}}$ and is not transmitted.

The DP protection module is trained in an adversarial manner with the help of a discriminator, as illustrated in Fig. 3. The discriminator, which consists of two FC layers followed by a sigmoid output layer, is only used for the training of the DP protection module. Specifically, it is trained to distinguish between $\hat{\mathbf{Z}}^1_{\text{private}}$ and $\hat{\mathbf{Z}}^2_{\text{private}}$, while the DP protection module is optimized to generate outputs that the discriminator cannot reliably differentiate from genuinely DP-noised data.

The training objective for the discriminator $D(\cdot)$ is defined using the binary cross-entropy (BCE) loss:

$$\mathcal{L}_D = -\mathbb{E}\left[\log D(\hat{\mathbf{Z}}^1_{\text{private}})\right] - \mathbb{E}\left[\log(1 - D(\hat{\mathbf{Z}}^2_{\text{private}}))\right]. \tag{19}$$

In contrast, the DP protection module $G(\cdot)$ is trained by minimizing the following loss:

$$\mathcal{L}_G = \mathbb{E}\left[\log(1 - D(\hat{\mathbf{Z}}^2_{\text{private}}))\right], \tag{20}$$

which encourages the DP protection module to generate perturbed latent codes whose noise patterns are indistinguishable from genuine DP noise. During training, the discriminator and the DP protection module are updated alternately. Moreover, the DP protection module must balance the trade-off between preserving reconstruction quality for the legitimate user and generating noise patterns that effectively resemble genuine DP noise. This trade-off will be further elaborated in Subsection E.

*2) Network Architectures:* The network architecture of the NN-based DP protection module is shown in Fig. 4. This module first vectorizes the input private latent codes $\mathbf{Z}_{\text{private}} \in \mathbb{R}^{m \times 512}$ into a one-dimensional representation $\mathbf{Z}_{\text{private}} \in \mathbb{R}^{512m}$. The vectorized signal is then passed through the FC layer and reshaped to its original dimension to produce the output $\hat{\mathbf{Z}}^2_{\text{private}} \in \mathbb{R}^{m \times 512}$. The NN-based DP deprotection module adopts the same architecture as the DP protection module. It vectorizes the input $\mathbf{Y}^1_{\text{private}}/\mathbf{Y}^2_{\text{private}} \in \mathbb{R}^{m \times 512}$,

passes it through an FC layer, and then reshapes it back to the original dimension to obtain $\hat{\mathbf{Y}}^1_{\text{private}}/\hat{\mathbf{Y}}^2_{\text{private}} \in \mathbb{R}^{m \times 512}$.

### D. The Calculation of the Sensitivity

In the context of DP for image protection [32], the sensitivity $\Delta f$ quantifies the maximum difference between the latent codes of any two distinct images in the training dataset. Accordingly, in this paper, $\Delta f$ is formally defined as

$$\Delta f \doteq \sup_{\mathbf{I}_1, \mathbf{I}_2 \in \mathcal{D}} \|f_{\text{inv}}(\mathbf{I}_1) - f_{\text{inv}}(\mathbf{I}_2)\|_2, \tag{21}$$

where $\mathbf{I}_1$ and $\mathbf{I}_2$ are two different images sampled from the training dataset $\mathcal{D}$, and $f_{\text{inv}}$ represents the inversion process of the Semantic StyleGAN. However, directly calculating $\Delta f$ using (21) presents two significant challenges. First, calculating the element-wise differences between the latent codes of all image pairs across a large training dataset results in considerable computational overhead. Second, certain source images may contain abnormal features that lead to outliers in the latent space, potentially distorting the sensitivity calculation.

To overcome these challenges, we adopt a clipping strategy that retains 99% of the latent codes, effectively mitigating the effect of outliers and simultaneously reducing the computational cost of calculating $\Delta f$. This process involves the following steps: (1) First, transform each image in the training dataset into its corresponding latent codes. (2) Next, determine the 0.5% quantile $a$ and the 99.5% quantile $b$ of all the latent code elements across the dataset. (3) Finally, adjust any elements outside the range $[a, b]$ to either $a$ or $b$.

After clipping, the sensitivity $\Delta f$ is calculated as:

$$\Delta f = \|b\mathbb{I}_n - a\mathbb{I}_n\|_2 = \sqrt{(b-a)^2 \cdot n}, \tag{22}$$

where $\mathbb{I}_n$ represents an all-ones vector of length $n$, and $n$ is the total number of elements in the latent codes for a single image. This approach efficiently calculates $\Delta f$ while addressing the computational and statistical challenges associated with large datasets and outliers.

### E. Training Strategy

In this subsection, we describe the training strategy employed for our system under two different settings for the eavesdropper.

*1) Proposed System (Basic Eavesdropper):* In the basic eavesdropper setting, the eavesdropper is unaware that DP noise with learnable pattern has been added to the private latent codes. Only the legitimate network (DP protection / deprotection modules of the legitimate user) and the discriminator are trained. The loss function for training the discriminator can be expressed as:

$$\mathcal{L}_{(1)} = -\mathbb{E}\left[\log D(\hat{\mathbf{Z}}^1_{\text{private}})\right] - \mathbb{E}\left[\log(1 - D(\hat{\mathbf{Z}}^2_{\text{private}}))\right]. \tag{23}$$

The loss function for training the legitimate network can be expressed as:

$$\mathcal{L}_{(2)} = \text{MSE}(\mathbf{Z}, \mathbf{S}_1) + \lambda \cdot \mathbb{E}\left[\log(1 - D(\hat{\mathbf{Z}}^2_{\text{private}}))\right]. \tag{24}$$

The first term of $\mathcal{L}_{(2)}$ measures the image reconstruction performance of the legitimate user. The second term of $\mathcal{L}_{(2)}$ measures how closely the DP noise with learnable pattern resembles the genuine DP noise, with $\lambda$ serving as the trade-off hyperparameter to balance these objectives. The two loss functions, $\mathcal{L}_{(1)}$ and $\mathcal{L}_{(2)}$, are optimized alternately.

*2) Proposed System (Stronger Eavesdropper):* In the stronger eavesdropper setting, the eavesdropper is aware that the semantic representation has been perturbed using DP noise with learnable pattern. The eavesdropper also has access to the network architecture of the DP deprotection module but does not know which latent codes have been protected. Unlike existing secure SemCom methods, in this setting, the performance of the eavesdropper is not considered during the training of the legitimate network (DP protection / deprotection modules), since the eavesdropper typically does not cooperate with the legitimate user. To simulate a realistic adversarial scenario, we allow the eavesdropper to train her network after the legitimate network has been trained. Specifically, we adopt a two-stage training strategy.

*First Training Stage:* This training stage follows the same procedure as that of the basic eavesdropper setting.

*Second Training Stage:* In this stage, only the eavesdropper's network (DP deprotection module of the eavesdropper) is trained. The loss function for this stage is:

$$\mathcal{L} = \mathrm{MSE}\left(\mathbf{Z}, \mathbf{S}_2\right), \tag{25}$$

which measures the image reconstruction performance of the eavesdropper.

## IV. PERFORMANCE EVALUATION

### A. Experimental Settings

*1) Pre-trained Model and Dataset:* We employ the pre-trained Semantic StyleGAN model trained on the CelebAMask-HQ dataset [38], where each image resized to $512 \times 512$. For the proposed system, we utilize the CelebAMask-HQ dataset for both training and testing. This dataset consists of 30,000 $1024 \times 1024$ RGB face images, with the first 28,000 images used for training and the remaining 2,000 images used for testing. All images are resized to $512 \times 512$ to maintain consistency with the pre-trained model.

*2) Private Latent Codes:* The latent codes of each image have a dimension of $28 \times 512$, where 28 represents the number of latent codes. Among these, the first two latent codes are shared latent codes, while the remaining latent codes are categorized as shape and texture codes (local latent codes). Each latent code has a length of 512. Thus, $H = (28 - 2)/2 = 13$, where $H$ represents the number of groups of local latent codes. Shared latent codes are always private latent codes. Additionally, certain latent codes from the local latent codes are selected as private latent codes.

*3) Privacy Budget and Channel SNR:* The sensitivity value $\Delta f$ is 351.88, and the privacy budget $\epsilon$ takes values from the set $\{1, 5, 10, 30, 100, 200, 300, 500, 800, 2000\}$. In this paper, we assume that the eavesdropper experiences the same channel conditions as the legitimate user, which represents the most challenging scenario for the legitimate user. Therefore, $\sigma_1 =$
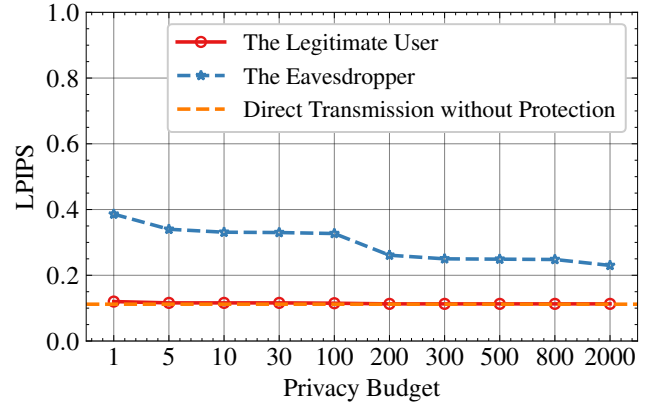


Fig. 5: The LPIPS performance of the *Proposed System (Basic Eavesdropper)* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.
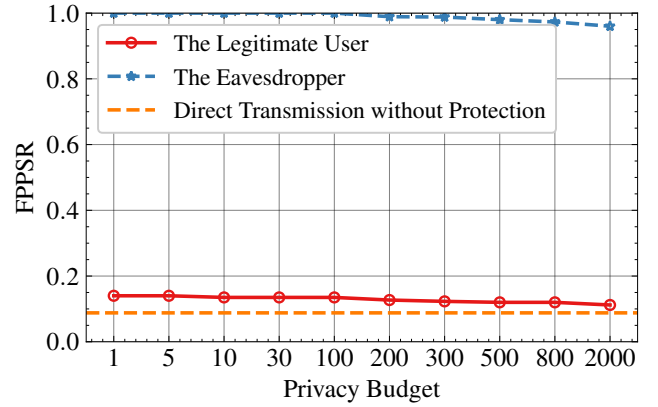


Fig. 6: The FPPSR performance of the *Proposed System (Basic Eavesdropper)* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.

$\sigma_2$, and $\mathrm{SNR} = \mathrm{SNR}_{\mathrm{leg}} = \mathrm{SNR}_{\mathrm{eve}}$. The channel SNR values are selected from the set $\{0, 5, 10, 15, 20\}$ dB.

*4) Training Settings:* Our experiments are conducted on a single NVIDIA RTX A6000 GPU. The batch size is set to 256, and the cosine annealing warm restarts optimizer is used to train the proposed system.

For the *Proposed System (Basic Eavesdropper)*, $\lambda$ is set to $1 \times 10^{-3}$. The shared latent codes, as well as the 4th to 7th latent codes, are selected as private latent codes. The initial learning rate is $3 \times 10^{-4}$, and the training lasts for 100 epochs.

For the *Proposed System (Stronger Eavesdropper)*, $\lambda$ is set to $1 \times 10^{-3}$. The shared latent codes and the 4th to 13th latent codes are private latent codes for the legitimate user. The eavesdropper guesses the shared latent codes and the 6th to 7th latent codes are private latent codes. In the first training stage, the initial learning rate is $3 \times 10^{-4}$, and the training lasts for 100 epochs. In the second training stage, the initial learning rate is $3 \times 10^{-4}$, and the training lasts for 50 epochs.

*5) The Benchmarks:* To evaluate the effectiveness of our proposed framework, we compare it against two representative benchmark methods.

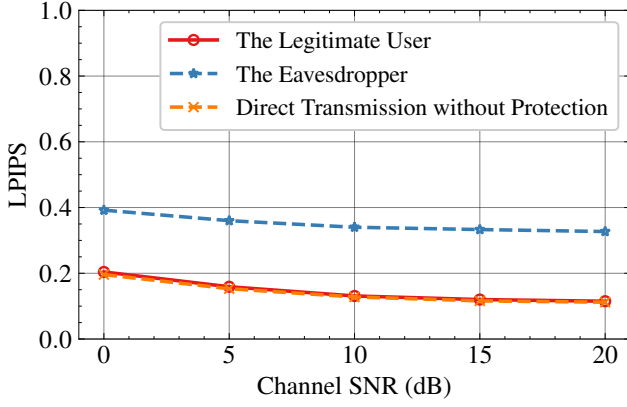In the first benchmark, we focus on scenarios without any security mechanisms. The latent codes $\mathbf{Z}$ are directly

Fig. 7: The LPIPS performance of the *Proposed System (Basic Eavesdropper)* under different channel SNRs, where the privacy budget $\epsilon$ is set to 100.
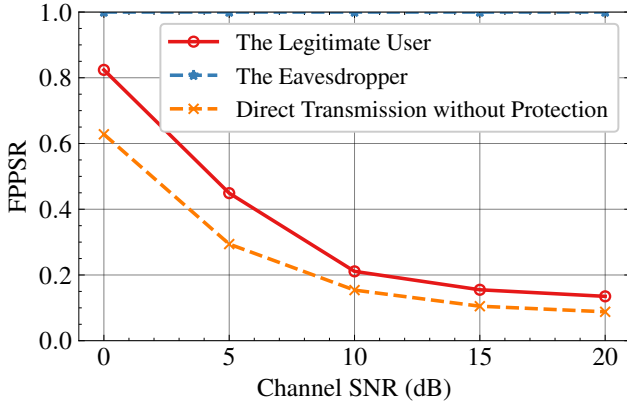


Fig. 8: The FPPSR performance of the *Proposed System (Basic Eavesdropper)* under different channel SNRs, where the privacy budget $\epsilon$ is set to 100.

transmitted from Alice to Bob without employing the DP protection module. Both Bob and Eve reconstruct the source image directly based on the received latent codes $\mathbf{Y}_1$ and $\mathbf{Y}_2$. Since they share the same channel SNR, their image reconstruction performances are identical. We refer to this benchmark as *Direct Transmission without Protection*. It is used to evaluate the effectiveness of privacy protection in our proposed system and its effect on the legitimate user's task performance.

The second benchmark, referred to as *Traditional DP Protection* [31], was proposed by Li *et al.* In this benchmark, we directly apply standard DP mechanisms [28], [31] to the private latent codes by adding Laplace noise calibrated according to a specified privacy budget. At the receiver, we train a dedicated NN that learns to remove the DP noise. This benchmark provides a formal privacy guarantee and serves as a baseline to evaluate the effectiveness of our proposed method in two aspects: its ability to mitigate the degradation in task performance caused by the non-invertibility of DP noise, and its capacity to provide effective image protection under DP constraints.

## B. Proposed System under the Basic Eavesdropper Setting

*1) Different Privacy Budgets:* Fig. 5 shows the LPIPS performance under different privacy budgets $\epsilon$, with the channel SNR fixed at 20 dB. The LPIPS metric quantifies perceptual similarity, where a lower value indicates better image reconstruction performance, and a higher value for the eavesdropper represents improved system security. Across the entire range of privacy budgets, our proposed method ensures that the legitimate user's LPIPS remains low and close to the benchmark value of 0.112, reflecting minimal degradation in reconstruction quality for the legitimate user. When $\epsilon = 1$, the LPIPS for the legitimate user is 0.120. As $\epsilon$ increases, the value decreases and stabilizes at 0.113 when $\epsilon \geq 200$, which is almost identical to *Direct Transmission without Protection*. In contrast, the LPIPS for the eavesdropper is 0.386 at $\epsilon = 1$, indicating severe reconstruction distortion. Although it gradually decreases as $\epsilon$ increases, it remains around 0.230 at $\epsilon = 2000$, which is considerably higher than that of the legitimate user. This highlights the ability of the DP protection module to degrade the reconstruction quality for the eavesdropper, thereby enhancing system security.

Fig. 6 shows the FPPSR performance under the same conditions. FPPSR reflects the system's ability to distinguish reconstructed face images from the original ones. Lower FPPSR values for the legitimate user indicate better image reconstruction quality, while higher FPPSR values for the eavesdropper imply stronger privacy protection. For the legitimate user, FPPSR begins at 0.140 when $\epsilon = 1$ and gradually declines to 0.112 at $\epsilon = 2000$. This value closely approaches the benchmark of 0.088. The eavesdropper consistently exhibits high FPPSR values across all privacy budgets. Starting at 1.0 when $\epsilon = 1$, the value only slightly decreases to 0.960 at $\epsilon = 2000$, indicating that the eavesdropper's reconstructed images are still highly distorted, unrecognizable, or slightly misleading. These results highlight the robustness of the *Proposed System (Basic Eavesdropper)* in achieving a strong balance between preserving reconstruction quality for the legitimate user and ensuring privacy protection against adversaries.

*2) Different Channel SNRs:* Fig. 7 shows the LPIPS performance over different channel SNRs, while the privacy budget $\epsilon$ is fixed at 100. We can see that our method achieves comparable reconstruction quality to *Direct Transmission without Protection* for the legitimate user across all channel SNRs, with LPIPS values only slightly higher than the benchmark. In contrast, the LPIPS values for the eavesdropper remain significantly higher than the benchmark at all channel SNRs, reflecting degraded perceptual similarity and thus enhanced system security. Specifically, at a channel SNR of 20 dB, the eavesdropper's LPIPS is 0.327, which is 0.212 higher than that of the legitimate user and 0.215 higher than the benchmark.

Fig. 8 shows the FPPSR performance under the same settings. Again, the legitimate user's FPPSR performance closely follows the benchmark and improves as SNR increases. Meanwhile, the eavesdropper's FPPSR remains consistently at 1.0 across all SNRs. This suggests that the eavesdropper consistently fails to reconstruct images that resemble the original identity, regardless of the channel quality. In conclusion,
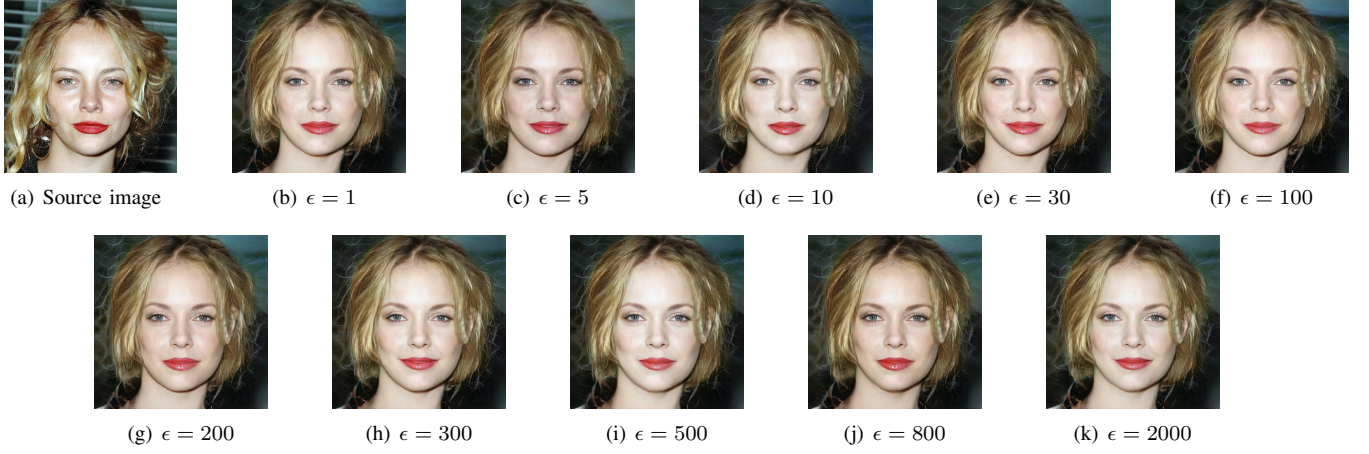
Fig. 9: Visual analysis of the reconstructed images by the legitimate user under different privacy budgets $\epsilon$ in the *Proposed System (Basic Eavesdropper)*. The channel SNR is set to 20 dB.
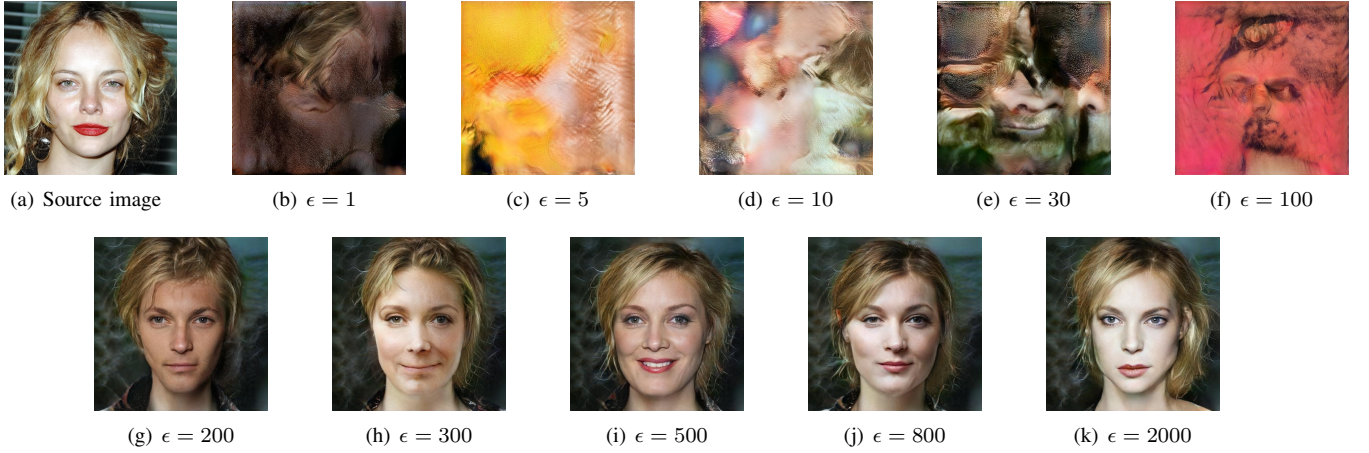


Fig. 10: Visual analysis of the reconstructed images by the eavesdropper under different privacy budgets $\epsilon$ in the *Proposed System (Basic Eavesdropper)*. The channel SNR is set to 20 dB.
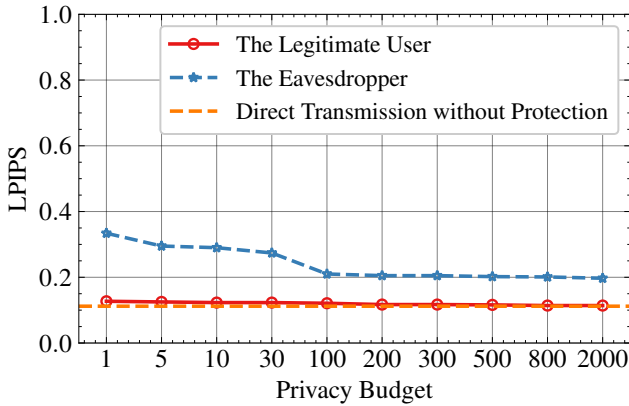


Fig. 11: The LPIPS performance of the *Proposed System (Stronger Eavesdropper)* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.
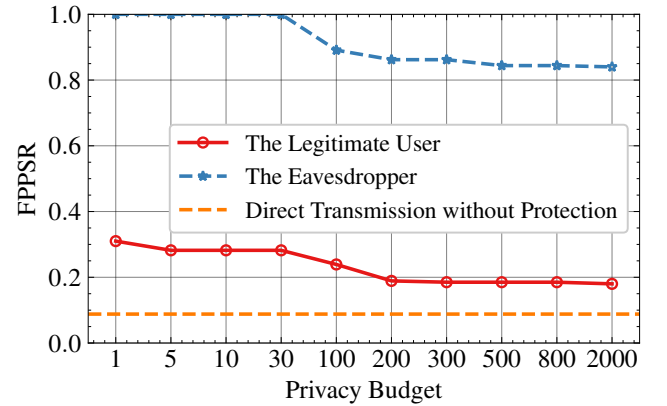
Fig. 12: The FPPSR performance of the *Proposed System (Stronger Eavesdropper)* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.

our method effectively maintains high task performance while substantially impairing the eavesdropper's ability to recover recognizable images, across all channel SNRs.

*3) Visual Evaluation:* Fig. 9 and Fig. 10 provide visual comparisons of the reconstructed images by the legitimate user and the eavesdropper under different privacy budgets $\epsilon$, with
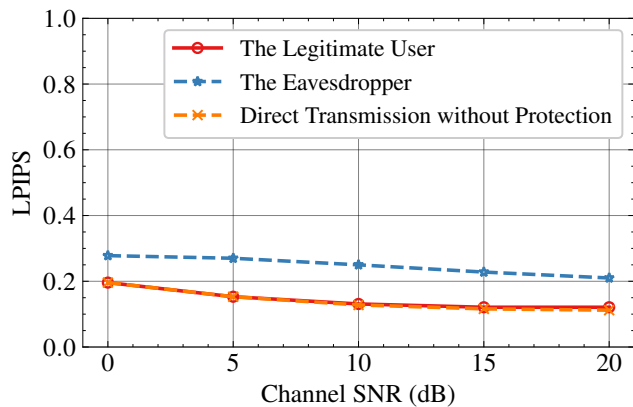
Fig. 13: The LPIPS performance of the *Proposed System (Stronger Eavesdropper)* under different channel SNRs, where the privacy budget $\epsilon$ is set to 100.
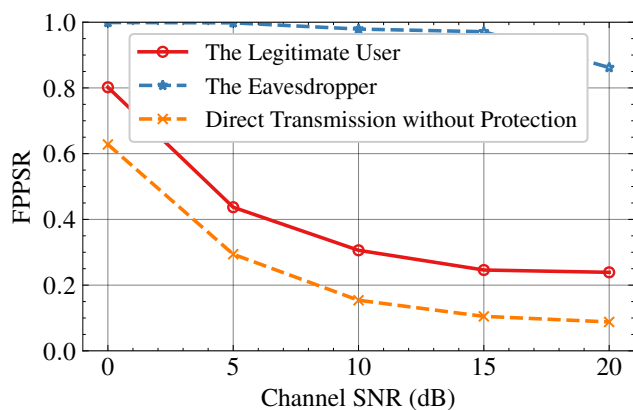


Fig. 14: The FPPSR performance of the *Proposed System (Stronger Eavesdropper)* under different channel SNRs, where the privacy budget $\epsilon$ is set to 100.

the channel SNR fixed at 20 dB. In Fig. 9, we observe the reconstruction results for the legitimate user across a range of privacy budgets. When $\epsilon = 1$, the reconstructed image already exhibits a high degree of semantic consistency with the source. As $\epsilon$ increases, the visual similarity further improves. When $\epsilon \geq 100$, the reconstructed images closely resemble the source image, demonstrating the robust deprotection capability of our DP deprotection module. These results indicate minimal perceptual differences and high-quality reconstruction.

In Fig. 10, we show the reconstructed results of the eavesdropper under the same range of privacy budgets. When $\epsilon$ is small (e.g., $\epsilon = 1, 5, 10, 30, 100$), the reconstructed images appear highly chaotic, lacking coherent facial structure, and are virtually unrecognizable. This indicates that our method effectively obfuscates sensitive semantic representation at low privacy budgets. As $\epsilon$ increases (e.g., $\epsilon = 200, 300, 500, 800, 2000$), the eavesdropper is able to recover structured face images. However, even at these higher privacy levels, the reconstructed images still exhibit significant deviations from the source image, making it visually implausible to infer the true identity. These results highlight the flexibility of the proposed approach, which can achieve varying levels of privacy protection by adjusting the privacy budget. By

carefully selecting the privacy budget, our proposed system can either effectively protect private information or generate fake yet artifact-mitigated images to slightly mislead the eavesdropper, ensuring robust system security while maintaining high reconstruction quality for the legitimate user.

### C. Proposed System under the Stronger Eavesdropper Setting

*1) Different Privacy Budgets:* Fig. 11 presents the LPIPS performance of the *Proposed System (Stronger Eavesdropper)* under different privacy budgets, with the channel SNR fixed at 20 dB. As shown in Fig. 11, the LPIPS values for the legitimate user remain consistently low across all privacy budgets, indicating strong reconstruction fidelity. Specifically, the LPIPS value begins at 0.127 when $\epsilon = 1$ and decreases to 0.114 as $\epsilon$ increases to 2000. These values are consistently close to the benchmark of 0.112, demonstrating that the NN-based DP deprotection module effectively mitigates the effect of the added DP noise with learnable pattern. In contrast, the LPIPS values for the eavesdropper remain higher than those for the legitimate user. When $\epsilon = 1$, the eavesdropper's LPIPS is 0.334, suggesting severely distorted reconstructions. Even as $\epsilon$ increases, the LPIPS decreases slowly and reaches 0.197 at $\epsilon = 2000$, which is still considerably higher than the legitimate user's LPIPS. This performance gap verifies the robustness of the proposed method in preserving semantic fidelity for the legitimate user while effectively preventing the eavesdropper from accessing private information.

Fig. 12 shows the corresponding FPPSR performance under varying $\epsilon$. For the legitimate user, the FPPSR decreases from 0.310 at $\epsilon = 1$ to 0.180 at $\epsilon = 2000$, gradually approaching the benchmark value of 0.088. This trend is consistent with the LPIPS results, indicating that the legitimate user can reconstruct high-quality images with increasing privacy budgets. On the other hand, the eavesdropper consistently exhibits very high FPPSR values, starting at 1.0 when $\epsilon = 1$ and decreasing to 0.840 at $\epsilon = 2000$. These results suggest that the eavesdropper's reconstructed images are either chaotic or substantially different from the original identity. The substantial FPPSR gap across all $\epsilon$ values further demonstrates the effectiveness of the proposed system.

*2) Different Channel SNRs:* Fig. 13 illustrates the LPIPS performance of the *Proposed System (Stronger Eavesdropper)* under different channel SNRs, with the privacy budget $\epsilon$ fixed at 100. As shown in Fig. 13, the LPIPS values for the legitimate user decrease from 0.196 at 0 dB to 0.121 at 20 dB, indicating a consistent improvement in perceptual reconstruction quality as the channel conditions improve. Notably, the LPIPS values closely approach the benchmark curve, demonstrating that the NN-based DP deprotection module effectively mitigates the effect of the added DP noise with learnable pattern and preserves high semantic fidelity for the legitimate user, even in low SNR regimes. In contrast, the eavesdropper consistently exhibits higher LPIPS values than the legitimate user across all SNR levels, starting at 0.278 when SNR is 0 dB and decreasing to 0.210 at 20 dB. Although the LPIPS values for the eavesdropper slightly decrease with increasing channel SNR, they consistently exceed those of the
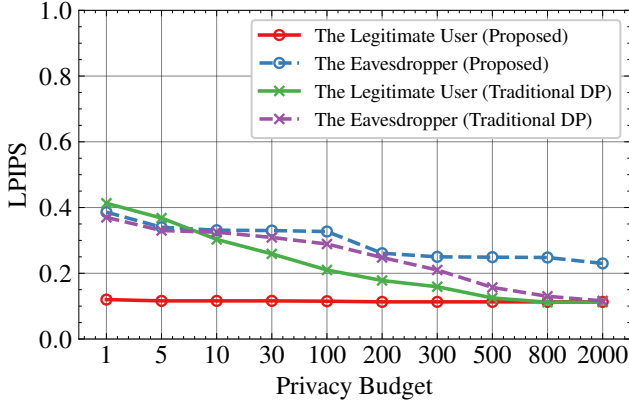
Fig. 15: The LPIPS performance of the *Proposed System (Basic Eavesdropper)* compared with the *Traditional DP Protection* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.
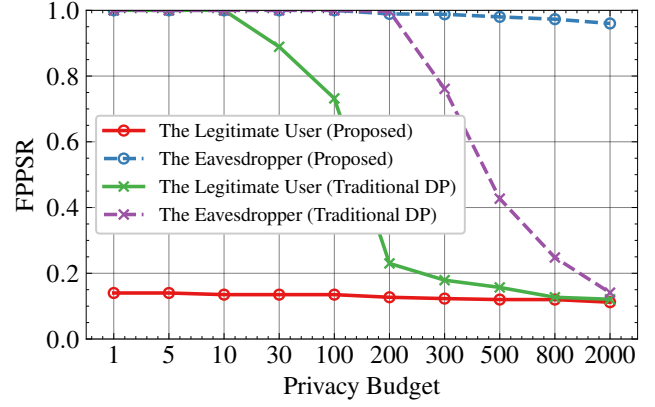


Fig. 16: The FPPSR performance of the *Proposed System (Basic Eavesdropper)* compared with the *Traditional DP Protection* under different privacy budgets $\epsilon$, where the channel SNR is set to 20 dB.

legitimate user by a large margin, confirming that the proposed method significantly degrades the perceptual similarity of the eavesdropper's reconstructed images. This performance gap highlights the proposed system's robustness in resisting unauthorized reconstruction, owing to the DP protection module and the knowledge gap, even when the eavesdropper has channel conditions equivalent to the legitimate user.

Fig. 14 presents the FPPSR performance under the same conditions as in Fig. 13. As shown in Fig. 14, the FPPSR for the legitimate user decreases significantly from 0.802 at 0 dB to 0.239 at 20 dB. This trend is consistent with the LPIPS results and further confirms that the preservation of facial identity improves as channel quality increases. Notably, the FPPSR values for the legitimate user remain close to the benchmark, reflecting strong identity retention capabilities. In contrast, the eavesdropper exhibits consistently high FPPSR values across all channel SNRs, beginning at 0.999 and reducing only to 0.862 at 20 dB. These elevated values indicate that the eavesdropper's reconstructed images are largely misidentified or appear chaotic, effectively obscuring the original facial identity. The persistent FPPSR gap further validates the effectiveness of the proposed system in protecting private information.

### D. Comparisons with Traditional DP Protection

In this subsection, we compare the performance of the *Proposed System (Basic Eavesdropper)* with the *Traditional DP Protection* benchmark [28], [31] under different privacy budgets $\epsilon$ in an equal-SNR (20 dB) wiretap channel scenario.

As shown in Figs. 15 and 16, our proposed method consistently outperforms the *Traditional DP Protection* benchmark in terms of both legitimate user performance and privacy protection against the eavesdropper. Specifically, the *Traditional DP Protection* benchmark introduces genuine DP noise to the private latent codes but suffers from inherent non-invertibility. As a result, the legitimate user's reconstruction performance deteriorates significantly, approaching the eavesdropper's performance. At low privacy budgets, the legitimate user struggles to reconstruct high-fidelity images. Conversely, at high privacy

budgets, the eavesdropper can still reconstruct images with high quality. More concretely, at low privacy budgets, both the legitimate user and the eavesdropper in the *Traditional DP Protection* benchmark show similarly poor LPIPS and FPPSR performance. At high privacy budgets, however, the eavesdropper achieves good LPIPS and FPPSR results, indicating effective image reconstruction and low system security. This shows that the *Traditional DP Protection* benchmark fails to provide effective privacy protection and high legitimate user performance in wiretap channel scenarios.

In contrast, our proposed method introduces the DP noise with learnable pattern, allowing the legitimate user to effectively mitigate its effect through the DP deprotection module. The results demonstrate that our method achieves significant improvements in both system security and task performance compared with the *Traditional DP Protection* benchmark.

## V. SUMMARY AND FUTURE WORK

This paper proposed a DP-based secure SemCom system over wiretap channels. Specifically, an NN-based DP protection module was introduced, where DP noise with learnable pattern was selectively added to private latent codes to provide privacy protection while maintaining high reconstruction quality for the legitimate user. To further enhance this process, a discriminator was employed to guide the noise generation toward resembling genuine DP noise, thereby enabling an approximate DP guarantee. At the receiver, a corresponding DP deprotection module was designed to effectively mitigate the effect of the introduced noise, enabling reliable image recovery for the legitimate user. Building on GAN inversion, the proposed fine-grained privacy protection strategy reduced computational overhead and model complexity compared with full-space DP protection. In addition, tunable privacy budgets provided flexible control over the system's security levels, which allowed the proposed system to produce either chaotic or slightly misleading images for the eavesdropper. Experimental results confirmed that our system outperformed both the previous DP-based method and direct transmission in terms

of system security and task performance, making it a practical and robust solution for secure SemCom.

This work can be extended in several promising research directions. First, to more effectively mislead potential eavesdroppers, future research could explore approaches for generating natural-looking protected images even under low privacy budgets. Second, incorporating robust anti-jamming strategies could help defend against semantic jamming attacks launched by malicious adversaries, thereby mitigating performance degradation for legitimate users.

## REFERENCES

[1] W. Chen, S. Tang, and Q. Yang, "Enhancing image privacy in semantic communication over wiretap channels leveraging differential privacy," in *Proc. 34th IEEE MLSP, London, UK, Sep. 2024*, pp. 1–6.

[2] P. Zhang, W. Xu, Y. Liu, X. Qin, K. Niu, S. Cui, G. Shi, Z. Qin, X. Xu, F. Wang, Y. Meng, C. Dong, J. Dai, Q. Yang, Y. Sun, D. Gao, H. Gao, S. Han, and X. Song, "Intellicise wireless networks from semantic communications: A survey, research issues, and challenges," *IEEE Commun. Surv. Tutorials*, vol. 27, no. 3, pp. 2051–2084, Jul. 2025.

[3] T. Han, Q. Yang, Z. Shi, S. He, and Z. Zhang, "Semantic-aware speech to text transmission with redundancy removal," in *Proc. IEEE ICC Workshops, Seoul, Korea, May 2022*, pp. 717–722.

[4] X. Peng, Z. Qin, X. Tao, J. Lu, and L. Hanzo, "A robust semantic text communication system," *IEEE Trans. Wirel. Commun.*, vol. 23, no. 9, pp. 11 372–11 385, Apr. 2024.

[5] T. Han, Q. Yang, Z. Shi, S. He, and Z. Zhang, "Semantic-preserved communication system for highly efficient speech transmission," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 1, pp. 245–259, Nov. 2022.

[6] X. Chen, J. Wang, L. Xu, J. Huang, and Z. Fei, "A perceptually motivated approach for low-complexity speech communication," *IEEE Internet Things J.*, vol. 11, no. 12, pp. 22 054–22 065, Mar. 2024.

[7] W. Chen, Y. Chen, Q. Yang, C. Huang, Q. Wang, and Z. Zhang, "Deep joint source-channel coding for wireless image transmission with entropy-aware adaptive rate control," in *Proc. IEEE GLOBECOM, Kuala Lumpur, Malaysia, Dec. 2023*, pp. 2239–2244.

[8] S. Tang, Q. Yang, L. Fan, X. Lei, A. Nallanathan, and G. K. Karagiannidis, "Contrastive learning-based semantic communications," *IEEE Trans. Commun.*, vol. 72, no. 10, pp. 6328–6343, May 2024.

[9] S. Wang, J. Dai, Z. Liang, K. Niu, Z. Si, C. Dong, X. Qin, and P. Zhang, "Wireless deep video semantic transmission," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 1, pp. 214–229, Nov. 2022.

[10] J. Guo, W. Chen, Y. Sun, J. Xu, and B. Ai, "Videoqa-sc: Adaptive semantic communication for video question answering," *IEEE J. Sel. Areas Commun.*, vol. 43, no. 7, pp. 2462–2477, Apr. 2025.

[11] M. Shen, J. Wang, H. Du, D. Niyato, X. Tang, J. Kang, Y. Ding, and L. Zhu, "Secure semantic communications: Challenges, approaches, and opportunities," *IEEE Netw.*, vol. 38, no. 4, pp. 197–206, Oct. 2023.

[12] Y. Li, Z. Shi, H. Hu, Y. Fu, H. Wang, and H. Lei, "Secure semantic communications: From perspective of physical layer security," *IEEE Commun. Lett.*, vol. 28, no. 10, pp. 2243–2247, Sep. 2024.

[13] W. Wang, Z. Tian, C. Zhang, and S. Yu, "SCU: an efficient machine unlearning scheme for deep learning enabled semantic communications," *IEEE Trans. Inf. Forensics Secur.*, vol. 20, pp. 547–558, Dec. 2024.

[14] J. Yang, S. Shao, F. Zou, and Y. Wu, "Dictionary learning-enabled privacy preserving semantic communication system," *IEEE Trans. Inf. Forensics Secur.*, vol. 20, pp. 5356–5371, May 2025.

[15] Y. Rong, G. Nan, M. Zhang, S. Chen, S. Wang, X. Zhang, N. Ma, S. Gong, Z. Yang, Q. Cui, X. Tao, and T. Q. S. Quek, "Semantic entropy can simultaneously benefit transmission efficiency and channel security of wireless semantic communications," *IEEE Trans. Inf. Forensics Secur.*, vol. 20, pp. 2067–2082, Jan. 2025.

[16] T. Marchioro, N. Laurenti, and D. Gündüz, "Adversarial networks for secure wireless communications," in *Proc. IEEE ICASSP, Barcelona, Spain, May 2020*, pp. 8748–8752.

[17] J. Shi, Q. Zhang, W. Zeng, S. Li, and Z. Qin, "Secure transmission in wireless semantic communications with adversarial training," *IEEE Commun. Lett.*, vol. 29, no. 3, pp. 487–491, Jan. 2025.

[18] E. Erdemir, P. L. Dragotti, and D. Gündüz, "Privacy-aware communication over a wiretap channel with generative networks," in *Proc. IEEE ICASSP, Virtual and Singapore, May 2022*, pp. 2989–2993.

[19] M. Zhang, Y. Li, Z. Zhang, G. Zhu, and C. Zhong, "Wireless image transmission with semantic and security awareness," *IEEE Wirel. Commun. Lett.*, vol. 12, no. 8, pp. 1389–1393, May 2023.

[20] X. Luo, Z. Chen, M. Tao, and F. Yang, "Encrypted semantic communication using adversarial training for privacy preserving," *IEEE Commun. Lett.*, vol. 27, pp. 1486–1490, Apr. 2023.

[21] Q. Qin, Y. Rong, G. Nan, S. Wu, X. Zhang, Q. Cui, and X. Tao, "Securing semantic communications with physical-layer semantic encryption and obfuscation," in *Proc. IEEE ICC, Rome, Italy, May 2023*, pp. 5608–5613.

[22] T. Tung and D. Gündüz, "Deep joint source-channel and encryption coding: Secure semantic communications," in *Proc. IEEE ICC, Rome, Italy, May 2023*, pp. 5620–5625.

[23] R. Meng, D. Fan, H. Gao, Y. Yuan, B. Wang, X. Xu, M. Sun, C. Dong, X. Tao, P. Zhang, and D. Niyato, "Secure semantic communication with homomorphic encryption," *arXiv:2501.10182v1 [cs.CR]*, Jan. 2025.

[24] W. Chen, S. Shao, Q. Yang, Z. Zhang, and P. Zhang, "A superposition code-based semantic communication approach with quantifiable and controllable security," *IEEE Trans. Mob. Comput. Early Access*, pp. 1–18, Sep. 2025.

[25] W. Chen, Q. Yang, S. Shao, Z. Shi, J. Chen, and X. Shen, "Can knowledge improve security? a coding-enhanced jamming approach for semantic communication," *arXiv:2504.16960v4 [cs.IT]*, Sep. 2025.

[26] B. He, Z. Chen, F. Wang, S. Wang, Z. Qin, and T. Q. S. Quek, "Diffusion-enabled secure semantic communication against eavesdropping," *arXiv:2505.05018v1 [cs.IT]*, May 2025.

[27] A. D. Wyner, "The wire-tap channel," *Bell Syst. Tech. J.*, vol. 54, no. 8, pp. 1355–1387, Oct. 1975.

[28] C. Dwork, "Differential privacy," in *Proc. ICALP, Venice, Italy, Jul. 2006*, pp. 1–12.

[29] L. Fan, "Image pixelization with differential privacy," in *Proc. IFIP WG, Bergamo, Italy, Jul. 2018*, pp. 148–162.

[30] ——, "Practical image obfuscation with provable privacy," in *Proc. IEEE ICME, Shanghai, China, Jul. 2019*, pp. 784–789.

[31] T. Li and C. Clifton, "Differentially private imaging via latent space manipulation," *arXiv:2103.05472v2 [cs.CV]*, Mar. 2021.

[32] H. Xue, B. Liu, M. Ding, T. Zhu, D. Ye, L. Song, and W. Zhou, "Dp-image: Differential privacy for image data in feature space," *arXiv:2103.07073v2 [cs.CR]*, Mar. 2021.

[33] Y. Wen, B. Liu, M. Ding, R. Xie, and L. Song, "Identitydp: Differential private identification protection for face images," *Neurocomputing*, vol. 501, pp. 197–211, Aug. 2022.

[34] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proc. IEEE/CVF CVPR, Salt Lake City, UT, USA, Jun. 2018*, pp. 586–595.

[35] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. NIPS, Lake Tahoe, Nevada, United States, Dec. 2012*, pp. 1106–1114.

[36] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proc. IEEE/CVF CVPR, Long Beach, CA, USA, Jun. 2019*, pp. 4690–4699.

[37] Y. Shi, X. Yang, Y. Wan, and X. Shen, "Semanticstylegan: Learning compositional generative priors for controllable image synthesis and editing," in *Proc. IEEE/CVF CVPR, New Orleans, LA, USA, Jun. 2022*, pp. 11 244–11 254.

[38] C. Lee, Z. Liu, L. Wu, and P. Luo, "Maskgan: Towards diverse and interactive facial image manipulation," in *Proc. IEEE/CVF CVPR, Seattle, WA, USA, Jun. 2020*, pp. 5548–5557.