# Diffusion-based Task-oriented Semantic Communications with Model Inversion Attack

Xuesong Wang, Mo Li, Xingyan Shi, Zhaoqian Liu, and Shenghao Yang, *Member, IEEE*

*Abstract*—Semantic communication has emerged as a promising neural network-based system design for 6G networks. Task-oriented semantic communication is a novel paradigm whose core goal is to efficiently complete specific tasks by transmitting semantic information, optimizing communication efficiency and task performance. The key challenge lies in preserving privacy while maintaining task accuracy, as this scenario is susceptible to model inversion attacks. In such attacks, adversaries can restore or even reconstruct input data by analyzing and processing model outputs, owing to the neural network-based nature of the systems. In addition, traditional systems use image quality indicators (such as PSNR or SSIM) to assess attack severity, which may be inadequate for task-oriented semantic communication, since visual differences do not necessarily ensure semantic divergence. In this paper, we propose a diffusion-based semantic communication framework, named DiffSem, that optimizes semantic information reconstruction through a diffusion mechanism with self-referential label embedding to significantly improve task performance. Our model also compensates channel noise and adopt semantic information distortion to ensure the robustness of the system in various signal-to-noise ratio environments. To evaluate the attacker's effectiveness, we propose a new metric that better quantifies the semantic fidelity of estimations from the adversary. Experimental results based on this criterion show that on the MNIST dataset, DiffSem improves the classification accuracy by $10.03\%$, and maintain stable performance under dynamic channels. Our results further demonstrate that significant deviation exists between traditional image quality indicators and the leakage of task-relevant semantic information.

*Index Terms*—Semantic communications, diffusion model, privacy preserving, model inversion attack.

## I. INTRODUCTION

**D**EEP learning based semantic communication system has emerged as a promising solution to enhance information transmission efficiency in sixth-generation (6G) wireless networks. Unlike traditional communication paradigms that prioritize the symbol-level accurate transmission, semantic communication shifts the focus to conveying the intended meaning through symbols [1]. In such systems, the transmitter and receiver collaboratively design mappings between meanings and channel symbols, naturally aligning with the principles of joint source-channel coding (JSCC) [2]. This approach accommodates more complex data sources (e.g., images) and diverse distortion measures, such as classification loss [3] or task-oriented metrics [4]. As a result, semantic communication has been studied for various types of sources such as text [5], [6], visual data [7]–[9] and video [10], [11],

The authors are with the School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, Guangdong 518172, China (e-mail: shyang@cuhk.edu.cn).

and has also been applied in edge computing [12], [13] and autonomous driving [14], [15]. By adopting a "comprehend-first, transmit-later" approach, semantic communication significantly reduces transmitted data volume, minimizes bandwidth consumption, and improves communication efficiency. However, when semantic communication is tailored to specific tasks, receivers predominantly rely on task-relevant semantic information, rendering the transmission and reconstruction of complete semantic data redundant and inefficient.

*Task-oriented semantic communication* addresses this limitation by directly optimizing performance for specific tasks. It enables the transmitter to extract only task-relevant semantic information and allows the receiver to focus on enhancing task-specific outcomes [16]. Resource allocation in such systems is tightly linked to task goals, making them particularly advantageous for resource-constrained environments or tasks with strict real-time requirements [17]. Studies have demonstrated the efficiency and robustness of task-oriented semantic communication in various scenarios, including image retrieval [18] and federated learning [19], as well as considering resource-constrained devices [20] and MIMO channels [21].

Despite its advantages, task-oriented semantic communication introduces new security risks. Attackers may exploit transmitted semantic features or intermediate representations to infer the original input data through model inversion attacks, leading to potential privacy leakage [22]. While transmitting only task-relevant information offers some level of privacy protection, adversaries can still reconstruct raw data or extract sensitive information using model inversion techniques [23]. Traditional privacy protection methods face significant challenges in the context of semantic communication. For instance, differential privacy [24] adds perturbations to model inputs or outputs which degrade semantic communication performance [25], while federated learning protects privacy through parameter sharing, but it is vulnerable to gradient leakage attacks [26].

Adversarial training methods based on information bottleneck have been proposed to reduce the pixel-level quality of reconstructed images [27]–[29]. However, these methods often evaluate privacy protection using image quality metrics like peak signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM), which do not directly reflect the semantic fidelity of the reconstructed information. For example, images with low SSIM may still retain critical semantic features, while high PSNR does not necessarily correlate with improved task performance [30]. This underscores the need to unify evaluation metrics for both attackers and recipients, focusing on semantic fidelity [16], [31] rather than image quality.

In this paper, we introduce a diffusion-based semantic communication system with model inversion attack involved. Diffusion models, as a generative artificial intelligence technique, offer potential solutions to these challenges. They can generate high-quality data (e.g., images, audio) through a gradual denoising process, recovering the original data distribution from random noise via reverse diffusion [32], [33]. By incorporating Gaussian white noise in the forward process and gradually denoising during sampling, diffusion models have demonstrated adaptability to additive Gaussian noise channels. Recent studies have explored their integration with semantic communication for tasks like audio restoration [34] and image denoising [35]–[37], primarily relying on conditional diffusion models equipped with receiver-side guidance.

Diffusion models excel in recovering and reconstructing semantic information due to their progressive denoising mechanism and ability to mimic underlying data distributions. Moreover, they can enhance privacy protection by allowing the sender to regulate distortion levels, thereby constraining adversarial behavior. However, applying diffusion models to privacy-sensitive semantic communication introduces unique challenges. Conditional diffusion models often assume the availability of guidance, which may not always be feasible. For instance, the sender may lack the capability to provide guidance due to performance limitations, or the receiver may need to keep task-specific details confidential. Additionally, in classification tasks, the guidance itself often overlaps with the basis for decision-making, creating a conflict between providing guidance and completing the task. These constraints underline the need for innovative approaches to leverage diffusion models in privacy-aware semantic communication systems.

In this paper, we propose a framework, named DiffSem, that leverages diffusion processes to support privacy-sensitive task-oriented semantic communication. Inspired by the forward process of diffusion models, we employ an adaptive noise injection mechanism, named self-noising, to regulate the quality of signal from the transmitter. Diffusion-based U-Net is applied at the receiver to enhance task performance, with the option to incorporate self-referential label embedding for further improvements. Inspired by the concept of semantic fidelity at the receiver [16], we propose a novel evaluation metric to more accurately quantify the semantic information obtained by adversaries through model inversion attacks. By introducing a strong classifier trained to evaluate adversaries' reconstructed images, our framework provides a direct assessment of the semantic fidelity achieved by adversaries while highlighting the inadequacy of traditional image quality metrics such as PSNR and SSIM in evaluating attack severity. Experimental results on the MNIST and CIFAR10 datasets demonstrate that DiffSem significantly improves receiver classification accuracy and achieves robust performance gains even under dynamic channel conditions.

Moreover, our findings reveal that traditional adversarial training often compromises the semantic fidelity of legitimate receivers without significantly reducing the semantic information extracted by adversaries, despite a marked deterioration in the quality of adversarial outputs. In contrast, our system implements a non-adversarial privacy protection strategy that enhances and preserves the semantic fidelity of legitimate receivers while restricting the fidelity attainable by adversaries. By conceptualizing adversaries as syntactic recovery agents and legitimate receivers as task-related semantic recovery agents, we explore the nuanced differences between syntactic and semantic information recovery. The self-denoising module in our system further helps refine this distinction by dynamically adjusting system outputs to optimize privacy protection and task performance. This unified framework not only advances task-oriented semantic communication but also provides a novel perspective on balancing task optimization and privacy protection in semantic communication systems.

In general, our contribution can be concluded as follows.

- We implement a diffusion-based semantic communication framework, named DiffSem, which leverages the unique properties of diffusion processes to optimize semantic information reconstruction at the receiver, and incorporates self-referential label embedding to significantly improve task performance. This framework ensures high-quality semantic recovery while providing robust performance under dynamic channel conditions. Experimental results on the MNIST dataset demonstrate that DiffSem boosts the receiver's classification accuracy by $10.03\%$ under high distortion conditions, while ensuring robust performance.

- We develop a non-adversarial privacy protection strategy that focuses on maintaining and enhancing the semantic fidelity of legitimate receivers while strictly limiting the upper bound of semantic fidelity attainable by adversaries. Our approach balances privacy protection and task performance by carefully regulating the semantic information distribution of system outputs. Additionally, the system employs an adaptive noise injection mechanism to further distinguish between semantic and syntactic information recovery, offering a novel perspective on privacy protection in semantic communication.

- We propose a new evaluation metric to accurately quantify the semantic information extracted by adversaries through model inversion attacks. By introducing a strong classifier to evaluate the semantic fidelity of reconstructed outputs, we provide a more precise and task-relevant framework for assessing privacy risks in semantic communication systems. In addition, we highlight a significant discrepancy between traditional image quality metrics (e.g., PSNR and SSIM) and task-relevant semantic information, underscoring the necessity of adopting semantic fidelity as a more appropriate evaluation criterion.

This paper is organized as follows. Section II describes the task-oriented semantic architecture as well as the model inversion attack, and general evaluations of the system. Section III introduces the modules and the training algorithms in our system in detail. Section IV includes the experiment results and discussion. Section V concludes the paper.

## II. SYSTEM ARCHITECTURE

In this section, we will first introduce the architecture of our system, then we will discuss the model inversion attack that
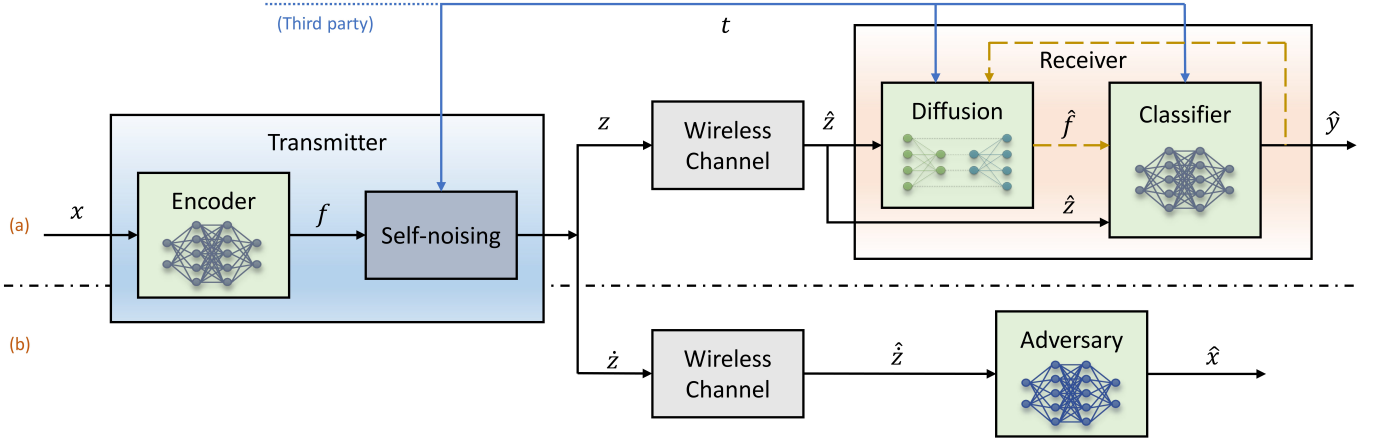
Fig. 1: System models. (a) Task-oriented transmission system. (b) Model inversion attack by adversary.

can estimate the input data. Next, we will discuss the general evaluations of privacy preservation in our system model.

### A. Diffusion based Task-oriented Semantic System

As shown in Fig. 1(a), assume there is a task-oriented transmitter-receiver pair that implements a specific task using the input raw data $x$ to get target object $y$. First, $x \in X$ is sent into a feature extractor module and a channel coding module to get neural-coded features $f \in F$. In our system, these two modules are both implemented by neural networks, and they can be regarded as a whole module named encoder $\Phi$ with parameter set $\phi$, like any other neural-network based joint-source-channel-coding (JSCC) system in related literature. Next, a self-noising module $\Psi$ is utilized to add a certain level of noise, based on system requirements, to get $z$ that is ready to transmit. These steps can be expressed as

$$z = \mathcal{T}_{(\phi)}(x) = \Psi(\Phi_\phi(x)) = \Psi(f) \quad (1)$$

where $\mathcal{T}$ means the transmitter. Note that the self-noising module $\Psi$ does not have any learnable parameters; details are introduced later.

After transmitting over the air, the receiver $\mathcal{R}$ obtains

$$\hat{z} = hz + \epsilon, \quad (2)$$

where $h$ is channel gain, $\epsilon$ is channel noise that has standard variation $\sigma$ in the wireless channel $\mathcal{C}$. Then, $\mathcal{R}$ performs its task in two steps. First, it uses a diffusion module $\Delta$ to perform denoising process that generates the estimated feature $\hat{f}$. Then $\hat{f}$ is channel decoded and is used to perform specific task. These two functions are also combined into one module called classifier $\Theta$, just analogous to the encoder module introduced above. The process of $\mathcal{R}$ can be defined as

$$\hat{y} = \mathcal{R}_{(\delta,\theta)}(\hat{z}) = \Theta_\theta(\Delta_\delta(\hat{z})) = \Theta_\theta(\hat{f}), \quad (3)$$

where $\delta$ and $\theta$ are the parameter set in denoising module $\Delta$ and classifier module $\Theta$, respectively, and $\hat{y}$ is the task output. Note that in wireless transmissions, $z, \hat{z} \in \mathbb{C}^\kappa$, where $\kappa$ means complex channel use. Since neural network only process real numbers, without further notice, let $k = 2\kappa$, and the length of $z$ ($f$ as well) is $k$ in our system's implementation.

There are two premises that need to be explained here. First, we assume that there is a trustworthy third party that can define the specific task and arrange models in each node of the system. Thus, we define that $\Theta$ is reachable only in $\mathcal{R}$ when models are trained and deployed by the third-party, which is a common prerequisite in task-oriented communications. Moreover, we argue that neither the plaintext nor the encrypted forms of $x$ and $y$ should be transmitted directly. Specifically, due to the limited computational capability of the transmitter, even computing $y$ locally is computationally challenging. As for $x$, we assume that both the receiver and potential eavesdroppers are curious and may attempt to infer information about $x$. Consequently, transmitting $x$—even in encrypted form—may still enable the receiver to recover $x$, unless extremely inefficient privacy-preserving inference mechanisms are employed. Therefore, to ensure data privacy, direct transmission of any form of $x$ or $y$ should be strictly avoided.

### B. Model inversion attack

Model inversion attack is an adversarial technique that targets machine learning models. The adversary can query the model and analyze how the model responds to attempt to infer the original input. It poses significant risks to data privacy, especially when sensitive information is involved.

As shown in Fig. 1(b), an adversary $\mathcal{A}$ may want to perform model inversion attack on the transmitter by recovering the estimation of $x$ (denoted by $\hat{x}$), given the detected $\hat{z}$ from a different wireless channel path. $\mathcal{A}$ tries to recover the estimation of $x$, denoted by

$$\hat{x} = \mathcal{A}_a(\hat{z}) = \mathcal{A}_a(\dot{h}z + \dot{\epsilon}), \quad (4)$$

where $\dot{h}$ and $\dot{\epsilon}$ are channel gain and noise in the adversary's channel, respectively, and $a$ is the parameter set of $\mathcal{A}$.

Assume $\mathcal{A}$ can access $\mathcal{T}$ repeatedly on the target device but cannot obtain the parameters and training data, and $\Theta$ is confidential to $\mathcal{A}$ and will not be accessed in any cases. Generally, the adversary first prepares a dataset that has a similar distribution as $X$, named $\dot{X}$, and sends $\dot{x}$ into $\mathcal{T}$ to
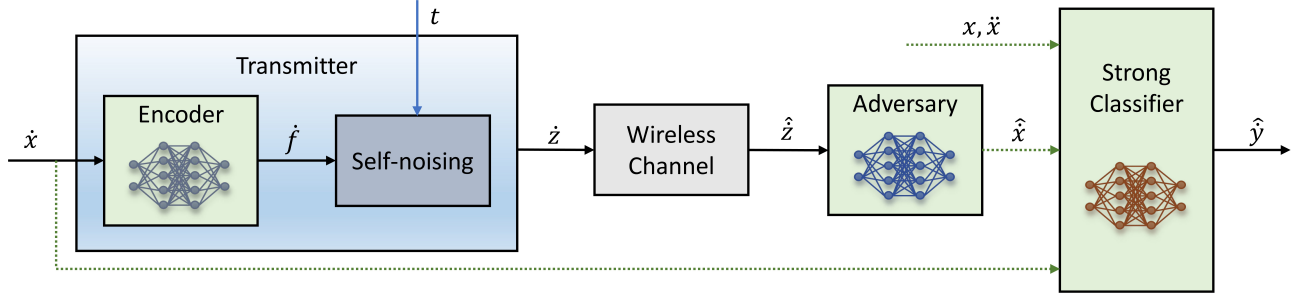
Fig. 2: Adversary prepares its dataset and accesses the transmitter repeatedly, with a well-trained strong classifier judging its outputs. $t$ is uniformly sampled from $[0, T]$. Green dash lines indicate the datasets that are used to train the strong classifier.

get $\hat{\hat{z}}$. When it gets enough data pairs $\{\dot{x}, \hat{\hat{z}}\}$, it trains on the data pairs using mean squared error (MSE) loss

$$\mathcal{L}_{\text{MSE}}(x_i, \hat{x}_i) = \frac{1}{M} \sum_{i}^{M} \|x_i - \hat{x}_i\|^2, \tag{5}$$

to minimize the reconstruction distortion. This is called a black-box setting [38].

### C. System Performance Evaluation

In this paper, we assume that the receiver performs the image classification task, so that $\mathcal{R}$ can be trained using cross-entropy (CE) loss

$$\mathcal{L}_{\text{CE}}(\phi, \theta) = -\sum_{i} p(y_i) log(q(\hat{y}_i)), \tag{6}$$

where $p(y)$ and $q(y)$ represent the true probability distribution and predicted probability of $y$, and the classification accuracy $\gamma$ can be defined as the number of correctly classified images divided by the total number of input images.

Generally, since $\mathcal{A}$ wants to recover the estimation of source images, PSNR and SSIM are commonly used in this scenario. PSNR is defined as

$$\text{PSNR} = 10 \log_{10} \left( \frac{R^2}{\text{MSE}} \right), \tag{7}$$

where $R$ is the maximum pixel value (e.g., 255 for 8-bit images), and MSE is the mean squared error. SSIM calculates the similarity between two images, expressed as

$$\text{SSIM}(x, \hat{x}) = \frac{(2\mu_x \mu_{\hat{x}} + c_1)(2\sigma_{x\hat{x}} + c_2)}{(\mu_x^2 + \mu_{\hat{x}}^2 + c_1)(\sigma_x^2 + \sigma_{\hat{x}}^2 + c_2)}, \tag{8}$$

where $\mu_x$ and $\mu_{\hat{x}}$ are the average pixel values of $x$ and $\hat{x}$, $\sigma_x^2$ and $\sigma_{\hat{x}}^2$ are the variances of $x$ and $\hat{x}$, $\sigma_{x\hat{x}}$ is the covariance between $x$ and $\hat{x}$, and $c_1$, $c_2$ are small constants that stabilize the division and avoid zero denominators. SSIM has a range of $[-1, 1]$ and usually observed within the range $[0, 1]$, where a larger value shows higher similarity between the two images.

In this paper, we give another index that measures the success rate of $\mathcal{A}$'s attack.

***Definition (1):*** *In task-oriented communication systems, the success rate of an attack is evaluated by the similarity between the task-relevant information in the attacker's estimated image and that in the original image.*

In our later experiments, we find that PSNR and SSIM may not be good indices in task-oriented system. PSNR and SSIM only indicate the similarity with original images in pixel-level. For a specific task, it is obvious that "task" information is more important than normal images, so focusing on how much task information (e.g., "label" in a category task) is recovered by $\mathcal{A}$ is more reliable, since we do not want $\mathcal{A}$ to acquire the purpose of the communication under the demand of privacy preserving. Details of implementation of ***Def. (1)*** are introduced in the next section.

### III. MODULES AND TRAINING

In this section, we will first introduce each module in our system model together with data flow, then we provide a detailed description of the training procedure based on our system design.

### A. Encoder, Classifier and Adversary Modules

In our system, encoder module $\Phi_\phi$ extracts neural-coded feature $f$ whose shape is $\sqrt{k} \times \sqrt{k}$ from source images $x^{CHW}$ with channel $C$, height $H$ and width $W$, and the compression rate is $\rho = k/(CHW)$. Feature $f$ is self-distorted to $z$ and is transmitted over the air, then the receiver gets channel-interfered signal $\hat{z}$. Receiver uses denoising module $\Delta_\delta$ and classifier module $\Theta_\theta$ to generate $\hat{y}$.

In related work [29], an auxiliary decoder $\mathcal{D}$ is commonly introduced to simulate an attacker during training, particularly in the context of adversarial learning frameworks. However, our approach does not adopt such attacker-aware training strategies. Instead, we define $\mathcal{A}$ as a model trained via model inversion attacks, solely for the purpose of evaluating the privacy risks after the transmission model is established. Importantly, $\mathcal{A}$ is not allowed to participate in the training or inference of the transmitter $\mathcal{T}$ and receiver $\mathcal{R}$. Likewise, if $\mathcal{D}$ is used in other methods, it should be restricted to the training phase and excluded from the evaluation phase. In our framework, since no adversarial training is involved, the auxiliary decoder $\mathcal{D}$ is entirely omitted.

To evaluate the semantic fidelity of adversary given ***Def. (1)***, we use a strong classifier $\bar{\Theta}$ that has more complex network structure than $\Theta$, to judge how much task-related information is contained in the image restored by the attacker $\mathcal{A}$, as shown in Fig. 2. It can be described by

$$\hat{y} = \bar{\Theta}_{\bar{\theta}}(\hat{x}), \tag{9}$$

where $\bar{\theta}$ is the parameter set of $\bar{\Theta}$. Based on ***Def. (1)***, we make a hard decision on $\bar{\Theta}$'s output to get $\hat{y}$ and compare it with $y$, which means we only select the category with the highest probability and ignore the soft information of other categories.

## B. Diffusion Module

We use denoising diffusion implicit models (DDIM) [39] to train the diffusion module in our system. In the forward process, diffusion module gradually adds Gaussian noise to the extracted feature $f$ with a series of schedules $\alpha_1, \ldots, \alpha_T \in [0, 1)$. To avoid ambiguity and facilitate discussion, in this section, let $f_0$ represent $f \in F$ mentioned earlier; that is, $f_0 \sim p(f_0)$ denotes the original feature, which is the output of encoder $\Phi$ and is not distorted yet. Then, noisy versions $f_1, \ldots, f_T$ are obtained by the following process in $T$ steps [32]:

$$p(f_t|f_{t-1}) = \mathcal{N}\left(f_t; \sqrt{1-\alpha_t}f_{t-1}, \alpha_t\mathbf{I}\right), \forall t \in \{1, \ldots, T\}, \tag{10}$$

where $\mathbf{I}$ is the identity matrix that has the same dimensions as $f_0$, and $\mathcal{N}(f_t; \mu, \sigma)$ is the normal distribution of mean $\mu$ and covariance $\sigma$ that generates $f_t$.

This recursive formulation allows $f_t$ to be sampled directly from

$$p(f_t|f_0) = \mathcal{N}\left(f_t; \sqrt{\bar{\alpha}_t}f_0, (1 - \bar{\alpha}_t\mathbf{I})\right), \forall t \sim \mathcal{U}\left(\{1, \ldots, T\}\right), \tag{11}$$

where $\mathcal{U}$ draws a normal distribution, $\bar{\alpha}_t = \prod_{i=1}^{t} \beta_i$ and $\beta_t = 1 - \alpha_t$. Then we have

$$f_t = \sqrt{\bar{\alpha}_t}f_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, \tag{12}$$

where $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$, to sample $f_t$ from $p(f_t|f_0)$ in the diffusion module's training step.

To ease the training and using of diffusion module, we choose $(\alpha_t)_{t=1}^{T}$ to be linear increasing constants, that is

$$\alpha_t = \alpha_1 + \left(\frac{\alpha_T - \alpha_1}{T - 1}\right)(t - 1), \tag{13}$$

where $\alpha_1 = 10^{-4}$, $\alpha_T = 2 \times 10^{-2}$. We train the diffusion module using $T = 500$, but $T' < T$ is used in the experiment to test the system's robustness. The values of $\sqrt{\bar{\alpha}_t}$ and $\sqrt{1 - \bar{\alpha}_t}$ varying by $t$ are shown in Fig. 3. We plot the "ESNR" line that is defined by

$$\mathcal{P} = \frac{\mathbb{E}(\|\sqrt{\bar{\alpha}_t}f_0\|^2)}{\mathbb{E}(\|\sqrt{1 - \bar{\alpha}_t}\epsilon_t\|^2)}, \tag{14}$$

where $\mathcal{P}$ represents the power ratio between the signal component $\sqrt{\bar{\alpha}_t}f_0$ and noise component $\sqrt{1 - \bar{\alpha}_t}\epsilon_t$ obtained after partially adding noise, and varies with the values of $\alpha_1$, $\alpha_T$ and $T$. In our system, $f_0$ is power-normalized after generated by encoder $\Phi$. Hence, according to Eq. (14), it is equivalent to the signal $f_0$ passing through an additive Gaussian white noise (AWGN) channel with signal-to-noise ratio $\mathcal{P}$, which we refer to as "equivalent SNR" (ESNR). As an example, we marked the point $t = 184$ in Fig. 3, where the energy of the signal component is approximately equal to that of the noise component, and the ESNR $p = 0.01$ dB $\simeq 0$ dB.

The signal $\hat{z}$ needs to be denoised first on the receiver side. Usually, a deep learning network called U-Net can be used to
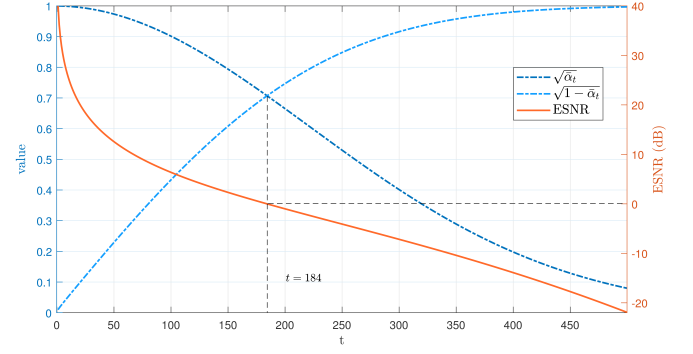


Fig. 3: Changes of $\sqrt{\bar{\alpha}_t}$, $\sqrt{1 - \bar{\alpha}_t}$ and $p$ with $t$.

predict $\epsilon_t$ during the sample process. U-Net is a convolutional neural network that can effectively handle multi-scale features and reconstruct high-fidelity output during the reverse process. U-Net learns to estimate $f_0$ from $f_t$ and $t$ by minimizing the difference between the predicted noise

$$\epsilon_{f_t} = \Delta_\delta(f_t, t) \tag{15}$$

and the actual noise $\epsilon_t$, and is trained using the MSE loss of these two noises, denoted by [32]

$$\mathcal{L}_{\text{U-Net}} = \mathbb{E}_{t \sim [1,T]}\mathbb{E}_{f_0 \sim p(f_0)}\mathbb{E}_{\epsilon_t \sim \mathcal{N}(0,\mathbf{I})} \|\epsilon_{f_t} - \epsilon_t\|^2. \tag{16}$$

With a well-trained U-Net, the goal is to progressively denoise $f_t$ and to recover $f_0$ starting from a latent variable $f_T \sim \mathcal{N}(0, \mathbf{I})$. Specifically, on step $t - 1$, given $f_t$ and $\epsilon_t$, we have

$$f_{t-1} = \sqrt{\bar{\alpha}_{t-1}}\left(\frac{f_t - \sqrt{1-\bar{\alpha}_t}\epsilon_{f_t}}{\sqrt{\bar{\alpha}_t}}\right) + \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_{f_t}. \tag{17}$$

To generate samples efficiently, DDIM allows to skip intermediate time steps. Instead of iterating over all $T$ steps, we can select a subset of time steps $\{t_1, t_2, \ldots, t_N\}$ with $N \ll T$. This sub-sampling strategy reduces the number of denoising steps while maintaining high-quality generation, making it significantly faster than traditional diffusion models.

There is another problem that needs to be discussed when applying diffusion to communication system. Generally, diffusion module cannot exactly recover the corresponding estimation of the input feature, as they are designed to learn the overall distribution of $f_0$ rather than the precise structure of individual inputs. When sampling from random white noise, the generation process is guided by learned distribution $f_0 \sim p(f_0)$, and U-Net produces outputs that have the same distribution as $p(f_0)$. But in communication systems, the transmitted information must be recovered in a specific order. To ensure the denoised signal meets this requirement, we can avoid adding excessive distortion to $f$ during the forward diffusion process, while proper guidance during the reverse process can further assist in generating outputs $\hat{f}$ that closely resemble the original $\hat{f}$.

Regarding the setting of $t$, let $0 \le T' \le T$, and perform

$$f_{T'} = \sqrt{\bar{\alpha}_{T'}}f_0 + \sqrt{1 - \bar{\alpha}_{T'}}\epsilon_{T'}, \tag{18}$$

in the self-noising module $\Psi$ before transmitting. In the receiver, $\Delta$ tries to recover $f_0$ using $T'$ steps' iteration that is the same as Eq. (17). This gives a way to control the weight

of "task-relevant information" $\sqrt{\bar{\alpha}_{T'}} f_0$ in the distorted output. The larger $T'$ is, the less "task-relevant information" in $f_{T'}$, and nearly no guidance involved when $T' = T$ (in which diffusion module generates $f_0$ without specific order). Since Eq. (18) is valid for any $0 \leq T' \leq T$, we can train the diffusion module to denoise from $T$ and deploy the well-trained module that is set by $T'$ when transmitting the signal, while receiver recovers the estimation of $f_0$ without other explicit guidance required. In this way, the choice of $T'$ controls the distortion of the transmitter, hence affects the difficulty for adversaries to estimate the original images with task information in them.

To further enhance system performance, the receiver can utilize self-generated label guidance. Specifically, the classifier module $\Theta$ generates the output $\hat{y}$ bypassing the diffusion module $\Delta$ at first. The accuracy of $\hat{y}$ is higher when distortion on $f$ is lower or channel condition is favorable, which means that $\hat{y}$ could be more helpful in such cases. Then $\hat{y}$ can be used to guide the diffusion module $\Delta$ in the denoising process (Fig. 1). Different from calculating the accuracy of $\hat{y}$, here we use the soft version of $\hat{y}$ that all the category info in the vector will be utilized as guidance. Next, $\hat{z}$ and $\hat{y}$ are fed into $\Delta$ to produce $\hat{y}'$, which is further reintroduced into $\Delta$ for refined denoising. By iterating this feedback loop multiple times, the accuracy of $\mathcal{R}$ is significantly improved. Consistent with guidance-based DDIM implementation, we provide explicit external control to $\Delta_\delta$ by modifying the noise estimation in Eq. (15) to

$$\epsilon_{(f_t, \hat{y})} = \Delta_\delta(f_t, \hat{y}, t), \quad (19)$$

where $\hat{y}$ serves as self-referential guidance. During the training, a classifier-free DDIM is trained with ground-truth labels $y$ by replacing the term $\epsilon_{f_t}$ with $\epsilon_{(f_t, y)}$ in Eq. (16). During sampling, $\hat{y}$ replaces $y$ to adaptively guide denoising. This soft label-embedding method ensures that:

- If $\hat{y}$ is more accurate, the denoising effect is enhanced;
- If $\hat{y}$ is unreliable (e.g., due to excessive self-noising in $\Psi$ or poor channel conditions), the system performance remains stable without degradation.

Although the external guidance from $\hat{y}$ may be suboptimal in extreme cases, it improves robustness in most scenarios. Detailed experimental validation of this iterative denoising framework is provided in Sec. IV.

### C. Diffusion Module Under Wireless Channel

Here we discuss the impact of wireless channel on the diffusion module in our system. To facilitate discussion, let $f_{T'} = z'$ and $h = 1$, then receiver gets

$$\hat{z}' = z' + \epsilon_{ch} = \left( \sqrt{\bar{\alpha}_{T'}} f_0 + \sqrt{1 - \bar{\alpha}_{T'}} \epsilon_{T'} \right) + \epsilon_{ch}, \quad (20)$$

where $\epsilon_{ch} \sim \mathcal{N}(0, \sigma_{ch}^2 \mathbf{I})$ represents channel noise. Obviously, the channel noise term adds extra distortion to $f_{T'}$ and changes its distribution, thus affects the sampling process at the receiver. To tackle this issue, we introduce some noise margin for possible channel interference during the forward process.

Suppose we need to distort $f_0$ to $f_{T'}$ which will be transmitted at $\mathcal{T}$, given specific timestep $T'$. Let $0 \leq T'' < T'$

and

$$f_{T''} = \sqrt{\bar{\alpha}_{T'}} f_0 + \sqrt{1 - \bar{\alpha}_{T''}} \epsilon_{T''}, \quad (21)$$

where $\epsilon_{T''} \sim \mathcal{N}(0, \mathbf{I})$. Note that the term $\sqrt{\bar{\alpha}_{T'}} f_0$ remains the same as $f_{T''}$. For simplicity, let $\epsilon_{T'} = \epsilon_{T''} = \epsilon$. Then, $\mathcal{T}$ transmits $f_{T''}$ to the wireless channel, and receiver obtains

$$\hat{z}'' = z'' + \epsilon_{ch} = \sqrt{\bar{\alpha}_{T'}} f_0 + \sqrt{1 - \bar{\alpha}_{T''}} \epsilon + \epsilon_{ch}. \quad (22)$$

If

$$\mathbb{E}\left( \|\sqrt{1 - \bar{\alpha}_{T'}} \epsilon\|^2 \right) = \mathbb{E}\left( \|\sqrt{1 - \bar{\alpha}_{T''}} \epsilon + \epsilon_{ch}\|^2 \right), \quad (23)$$

then $\mathbb{E}\left( \|\hat{z}''\|^2 \right) = \mathbb{E}\left( \|z'\|^2 \right)$, which is equivalent to that $\mathcal{R}$ receives $f_{T'}$. So we can transmit $z''$ while maintaining the distribution after the wireless transmission, and then the diffusion module can recover the estimation of $f_0$ without further ambiguity.

Now we need to find $T''$. Let $\gamma$ represent channel's signal-noise power ratio, then we have

$$\sigma_{ch}^2 = \frac{\mathbb{E}\left( \|z''\|^2 \right)}{\gamma} = \frac{\bar{\alpha}_{T'} \mathbb{E}\left( \|f_0\|^2 \right) + (1 - \bar{\alpha}_{T''})D}{\gamma}, \quad (24)$$

where $D$ is the length of $\epsilon$ and $\mathbb{E}[\|\epsilon^2\|] = D$.

By considering Eq. (22) to (24), we can derive that

$$\bar{\alpha}_{T''} = \frac{\bar{\alpha}_{T'} \left( \gamma + \mathbb{E}\left( \|f_0\|^2 \right) \right) + D}{\gamma + D}. \quad (25)$$

Details can be found in Appendix A. In Eq. (25), $\bar{\alpha}_t$ can be regarded as discrete function of $t$ according to Eq. (13), so $T''$ can be approximated given $T'$, $\gamma$ and $f_0$ [1]. In this way, our diffusion module can effectively handle the influences from the wireless channel [2].

If $h$ is involved (e.g., Rayleigh channel), then

$$\hat{z}'' = h\sqrt{\bar{\alpha}_{T'}} f_0 + h\sqrt{1 - \bar{\alpha}_{T''}} \epsilon + \epsilon_{ch}. \quad (26)$$

Let

$$\epsilon_{total} = h\left( \sqrt{1 - \bar{\alpha}_{T''}} \epsilon_{T''} - \sqrt{1 - \bar{\alpha}_{T'}} \epsilon_{T'} \right) + \epsilon_{ch}, \quad (27)$$

and we can get

$$\frac{\hat{z}''}{\hat{h}} = f_{T'} + \frac{\epsilon_{total}}{\hat{h}} \quad (28)$$

Receiver can perform the signal estimation using minimum mean square error (MMSE) once $\hat{h}$ is estimated. However, the discussion of $h$ here is beyond the scope of our paper. In this paper, we focus on improving the task execution performance when model inversion attack occurs due to neural-based nature of the semantic systems, and we use AWGN channel as an example to illustrate our points. We will leave the study of systems involving more complex channel conditions in the future work.

### D. System Training Procedure

Here we introduce the training procedure in our experiment. As shown in Fig. 1, $t$ is used as a global parameter that controls

---

[1] $f_0$ and $\epsilon_{ch}$ are of the same length.

[2] In the implementation, we find the maximum $T''$ that $\alpha_{T''}$ is no larger than the right side of Eq. (25). If cannot, let $T'' = 0$; generally, it means that $T'$ is small enough, or $\gamma$ is larger enough. Both situations mean that there will not be significant distortion on $f_0$.

---

**Algorithm 1** System Training Preparation

---

1: **Initialization:** Initialize weights $W$ and biases $b$ in $\mathcal{T}$, $\mathcal{R}$, $\mathcal{A}$ and $\bar{\Theta}$.
2: **Input:** $x \in X$ and $\dot{x} \in \dot{X}$.
3: **Stage 1: Training $\Phi_\phi$ and $\Theta_\theta$.**
4:     $\Phi_\phi(x) = z$.
5:     Transmit $z$ over wireless channel and get $\hat{z}$.
6:     $\Theta_\theta(z) = \hat{y}$.
7:     Calculate $\mathcal{L}_{\text{CE}}$ by Eq. (6) and update $\phi$, $\theta$.
8: **Stage 2: Training $\Delta_\delta$**
9:     $\Phi_\phi(x) = f$.
10:     Train U-Net by minimizing $\mathcal{L}_{\text{U-Net}}$ using Eq. (16), and update $\delta$.
11: **Stage 3: Training $\mathcal{A}_a$ and $\bar{\Theta}_{\bar{\theta}}$**
12:     $\Phi_\phi(\dot{x}) = \dot{z}$.
13:     Transmit $\dot{z}$ over wireless channel and get $\hat{\dot{z}}$.
14:     $\mathcal{A}_a(\hat{\dot{z}}) = \hat{x}$.
15:     Directly distort $x$ to get $\ddot{x}$.
16:     $\bar{\Theta}_{\bar{\theta}}(cat(x, \dot{x}, \hat{x}, \ddot{x})) = \hat{\dot{y}}$.
17:     Use $\hat{x}$ to calculate MSE loss by Eq. (5); use $\hat{\dot{y}}$ to calculate CE loss by Eq. (6).
18:     Update $a$ and $\bar{\theta}$.
19: **Output:** Pre-trained network $\mathcal{T}$, $\mathcal{R}$, $\mathcal{A}$ and $\bar{\Theta}$.

---

**Algorithm 2** Receiver Fine-tuning

---

1: **Initialization:** Load pre-trained networks $\mathcal{T}$, $\mathcal{R}$ and $\mathcal{A}$.
2: **Input:** $x \in X$.
3: **Given** $0 \le T' \le T$,
4:     Find $T''$ by Eq. (25).
5:     $\mathcal{T}_\phi(x, T', T'') = z$.
6:     Transmit $z$ over wireless channel and get $\hat{z}$.
7:     $\Theta_\theta(cat(\hat{z}, T')) = \hat{y}$.
8:     **if** using E-DiffSem:
9:       **while** accuracy of $\hat{y}$ increases:
10:         **for** $t \in \{t_1 = 0 < t_2 < \cdots < t_N = T'\}$:
11:           Denoise $\hat{z}$ by Eq. (17) and Eq. (19) from $f_t$ to $f_{t-1}$, with $\hat{y}$ involved.
12:         **end for** and get $\hat{f}$.
13:         $\Theta_\theta(cat(\hat{f}, \hat{z}, T')) = \hat{y}$.
14:       **end while**
15:     **elif** using DiffSem:
16:       **for** $t \in \{t_1 = 0 < t_2 < \cdots < t_N = T'\}$:
17:         Denoise $\hat{z}$ by Eq. (17) and Eq. (15) from $f_t$ to $f_{t-1}$.
18:       **end for** and get $\hat{f}$.
19:       $\Theta_\theta(cat(\hat{f}, \hat{z}, T')) = \hat{y}$.
20:     **end if**
21:     Calculate $\mathcal{L}_{\text{CE}}(y, \hat{y})$ by Eq. 6 and update $\theta$.
22: Change $T'$ and return to Line 3 to repeat.
23: **Output:** Finetuned network $\Theta_\theta$.

---

the distortion at the transmitter and helps the denoising at the receiver. All modules are neural-based except for self-noising module and wireless channel modules. As for wireless channel, we consider AWGN channel in the experiment. The self-noising module can be regarded as a channel compensator, ensuring that the signals received by the diffusion module appear as if they have passed through a Gaussian channel, thereby conforming to the input requirements for denoising. MNIST and CIFAR10 are used as datasets to train the system model. Given the small size of images and features in these datasets, we only employ a simple diffusion module that is easy to train and implement, as well as convenient for deployment considering the low-latency requirements of practical communication systems.

For ease of tracking, we divide the training procedure into two parts. The first part consists of three stages, each designed to prepare a specific module, as outlined in Algorithm 1. We assume that when an adversary $\mathcal{A}$ attempts to "join" a transmitter-receiver pair, the transmitter $\mathcal{T}$ and receiver $\mathcal{R}$ have already been fully trained and are successfully deployed on the transmission nodes. Consequently, the training of $\mathcal{A}$ must occur after the completion of the $\mathcal{T}$-$\mathcal{R}$ pair. As shown in Fig. 2, we use the well-trained transmitter to train the adversary, where receiver is not involved here. Hence, the performance of $\mathcal{A}$ will inherently depend on the training outcomes of $\mathcal{T}$. By optimizing the training of the $\mathcal{T}$-$\mathcal{R}$ pair and subsequently training $\mathcal{A}$ to its fullest potential, we aim to establish a (non-strict) upper bound on the amount of valid information the adversary can extract under the given dataset and system. The training of strong classifier $\bar{\Theta}$ can be carried out simultaneously with $\mathcal{A}$, or after the training of $\mathcal{A}$ is finished.

The second part is fine-tuning the receiver given fixed $\Phi$

and $\Psi$, as described in Algorithm 2. We implemented Diff-Sem, which performs a single-pass pipeline processing on the receiver's input, and an enhanced version named E-DiffSem, where "E" stands for "enhance". E-DiffSem utilizes soft $\hat{y}$ as self-referential guidance to improve the performance of denoising. Specifically, self-referential $\hat{y}$ is iteratively utilized until the accuracy of $\hat{y}$ reaches stability. For both methods, the diffusion steps $N$ are set to $\min(T', 50)$. Experimental results and discussions will be elaborated in the next section.

## IV. EXPERIMENTS AND RESULTS

In this section, we will first explain the module design in our system and introduce the details of experiment initialization. Next, we present the results and discussions based on our experimental findings.

### A. Module Implementation

The structures of modules involved in our system are as shown in Table I and Table II. We evaluate our proposed system on MNIST and CIFAR10 datasets. Functions of each module are as mentioned above. The compression rates ($\rho$) are set to $8.16\%$ and $6.25\%$ for the MNIST and CIFAR-10 datasets, respectively. These values are chosen merely for the convenience of model training, without any specific or intentional numerical design. Note that before the encoder's output is sent to the self-denoising module, it undergoes power normalization to make sure that $\sum_i(f_i^2/k) \le 1$, where $i$ is the pixel index. This operation is performed during the training of the $\mathcal{T}$-$\mathcal{R}$ pair and the training of $\Delta$, to facilitate the denoising

TABLE I
SUMMARY OF MODULES FOR MNIST DATASET

| Modules | Layers | Output Dimensions |
|---|---|---|
| **Encoder** | Conv2d+LeakyReLU | $16\times14\times14$ |
| | Conv2d+LeakyReLU | $32\times8\times8$ |
| | Conv2d+LayerNorm | $1\times8\times8$ |
| **Classifier (in Receiver)** | Linear+ReLU | $64\times2+1$ |
| | Linear | 10 |
| **Strong Classifier (in Third-party)** | Linear+Tanh | 784 |
| | Linear+Tanh | 128 |
| | Linear+Sigmoid | 10 |
| **Adversary** | Linear+Tanh | 784 |
| | Linear | $1\times28\times28$ |

TABLE II
SUMMARY OF MODULES FOR CIFAR10 DATASET

| Modules | Layers | Output Dimensions |
|---|---|---|
| **Encoder** | Conv2d+LeakyReLU | $16\times32\times32$ |
| | Conv2d+LeakyReLU | $32\times16\times16$ |
| | Residual | $32\times16\times16$ |
| | Conv2d+LeakyReLU +LayerNorm | $3\times8\times8$ |
| **Classifier (in Receiver)** | Conv2d+Tanh | $64\times8\times8$ |
| | Conv2d+Tanh | $512\times4\times4$ |
| | Residual Block | $512\times4\times4$ |
| | AdaptiveAvgPool2d | 512 |
| | Linear+Tanh | 10 |
| **Strong Classifier (in Third-party)** | Conv2d+ReLU | $64\times32\times32$ |
| | Conv2d+ReLU | $256\times16\times16$ |
| | Residual Block | $512\times8\times8$ |
| | AdaptiveAvgPool2d | $512\times2\times2$ |
| | Linear+ReLU | 256 |
| | Linear | 10 |
| **Adversary** | Conv2d+Tanh | $512\times8\times8$ |
| | (ConvTrans2d+Tanh +Residual Block)$\times2$ | $128\times32\times32$ |
| | Conv2d+Tanh | $3\times32\times32$ |

process at the receiver. Therefore, the output of the $\mathcal{T}$ does not require further normalization.

There are two classifiers in the system, $\Theta$ and stronger $\bar{\Theta}$. When using the MNIST dataset, the input of $\Theta$ is a combination of $\hat{z}$, $\hat{f}$ and $t$, resulting in a dimension of $64 \times 2 + 1$, while the input size of $\bar{\Theta}$ is $1 \times 28 \times 28$. For the CIFAR10 dataset, the input of $\Theta$ is a concatenation of $\hat{z}$ and $\hat{f}$, with $t$ padded alongside them. Hence, the dimension of the input becomes $(3 \times 2) \times 8 \times (8+1) = 6 \times 8 \times 9$. The strong classifier $\bar{\Theta}_{\bar{\theta}}$ needs to learn from the combined dataset that includes training dataset $X$, adversary's dataset $\dot{X}$ and the outputs of the adversary, and aims to achieve as accurate classification as possible. $\bar{\Theta}$ also learns to categorize images that are directly distorted according to channel conditions, denoted as $\ddot{x}$.

TABLE III
SUMMARY OF STRUCTURE FOR U-NET

| Processing Steps | Layers | Output Dimensions |
|---|---|---|
| **Encoding** | ResidualConv Block | $128\times8\times8$ |
| | Unet Down | $128\times4\times4$ |
| | Unet Down | $256\times2\times2$ |
| **Vectorization** | AvgPool2d+GELU | $768\times1\times1$ |
| **Embedding** | 2$\times$EmbedFC (Time) | $(256/128)\times1\times1$ |
| | 2$\times$EmbedFC (Context) | $(256/128)\times1\times1$ |
| **Decoding** | ConvTrans2d +GroupNorm+ReLU | $256\times2\times2$ |
| | UnetUp | $128\times4\times4$ |
| | UnetUp | $128\times8\times8$ |
| **Output** | Conv2d+GroupNorm +ReLU+Conv2d | $1\times8\times8$ |

The structure of U-Net in the diffusion module is shown in Table III [3]. We employ U-Net to denoise the input effectively while incorporating the guide information. Specifically, it first processes noisy features by down-sampling them through an encoding step. The down-sampled features are then vectorized and combined with time embeddings and context embeddings (if any) that guide the decoding step. Decoding step up-samples the features and integrates skip connections from the encoding. Finally, convolutional layers reconstruct the denoised features that match the input dimensions. This structure allows the model to effectively denoise images while incorporating time and context information. In our implementation, the diffusion module remains static after training, allowing pre-trained diffusion models to handle and generate transmission information efficiently. This eliminates the need for joint training, significantly lowering both computational overhead and energy usage [37].

Other relevant details are outlined as follows. We use maximum diffusion timestep $T = 500$ and $0 \le T' \le T$. We set the learning rate varying from $1 \times 10^{-3}$ to $1 \times 10^{-4}$. Adam optimizer was utilized for training, along with the ReduceL-ROnPlateau scheduler to dynamically adjust the learning rate, ensuring efficient convergence and preventing overfitting. Both MNIST and CIFAR10 datasets are divided into two portions, with $60\%$ allocated for the $\mathcal{T}$-$\mathcal{R}$ pair and $40\%$ designated for $\mathcal{A}$.

We reproduced several baseline models to provide a clear comparison and better evaluate the performance of our system.

- Information bottleneck theory and adversarial learning (named "IBAL") [29], that could extract and transmit task-related features while enhancing privacy by significantly decreasing the SSIM and PSNR values of estimated images that the adversary could generate;
- IBAL under dynamic channel conditions (named "IBAL-D") [29], that can adjust the weights of the classifier's decision term and the attacker's recovery term in the loss

---

[3] Due to limited space and to avoid excessive repetition, here we present the structure of U-Net used in E-DiffSem. For DiffSem, context embedding is not involved, the vectorization step generates $128\times2\times2$, and the output dimensions with 256 should be 128.
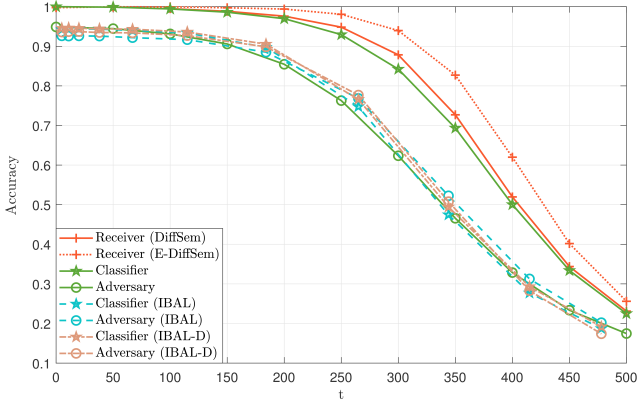
Fig. 4: System performances and baselines using MNIST dataset.
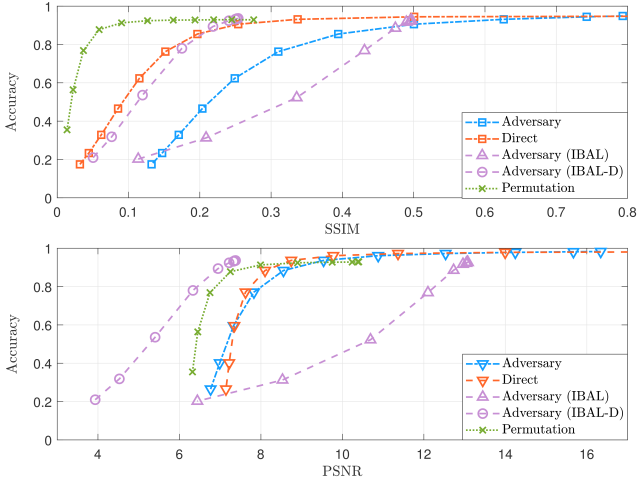


Fig. 5: Plot of accuracy versus SSIM and PSNR using MNIST dataset.

function based on the channel conditions.

- Direct mode (named "Direct"), where images $x \in X$ are directly distorted according to the channel conditions, and then are classified by the strong classifier.
- Permutation mode (named "Permutation" [40]), where the order of elements of images are rearranged according to a specific permutation rule or key, without altering the actual content of the elements. Then permuted images are classified by the strong classifier.

### B. Results on MNIST dataset

Fig. 4 provides a comprehensive performance comparison of various semantic communication systems for the MNIST dataset with a particular focus on accuracy, under channel SNR 15dB. Here we provide another curve named "Classifier", which means $\hat{z}$ is directly sent into the classifier $\Theta_\theta$ and bypasses the diffusion module. The green solid lines with notations stars and circles represent the classification accuracy of $\hat{z}$ and $\hat{x}$, without other modules involved. They can be simply regarded as a classification task and an image restoration task under different channel conditions that are determined by $t$, and they reveal a fundamental gap that reflects the intrinsic difference between recovering semantic information (task-specific features) and syntactic information (pixel-level details) given specific task and dataset, confirming

that task-oriented recovery is inherently more efficient. For example, at $t = 300$, the classifier achieves $84\%$ accuracy, while the adversary's reconstructed images yield only $62\%$ accuracy for the same task. Since semantic information is a highly condensed version of syntax and requires less feature to describe, reconstructing syntax is inherently more demanding than recovering semantics.

The two red lines at the top of the figure represent the classification accuracy of the receiver after incorporating the diffusion module. It can be observed that using DiffSem provides a slight improvement in the receiver's recovery, while the label-embedded E-DiffSem significantly enhances system performance. This aligns with the earlier discussion that $\sqrt{\bar{\alpha}_t} f_0$ can be regarded as an intrinsic guide, and the receiver's preliminary recovery $\hat{y}$ can serve as guidance to effectively boost diffusion. Most notably, when $t = 300$, DiffSem and E-DiffSem improve classification accuracy by $3.35\%$ and $10.03\%$ comparing to the system without diffusion, respectively. One possible explanation for the superior performance of diffusion-based methods lies in their ability to model complex distributions and generate high-quality reconstructions. Diffusion models iteratively refine the output by design, which aligns well with the requirements of semantic communication, where accurate reconstruction of semantic information is crucial. E-DiffSem further optimizes the diffusion process by effectively leveraging the label estimated by the classifier, leading to higher accuracy.

The IBAL and IBAL-D baselines provide an interesting contrast to the proposed methods. Note that the curves in Fig. 4 use SNR as the horizontal axis during training; here they have been recalculated based on Eq. (14) and presented with $t$ as the horizontal axis. We employ the same strong classifier structure for all adversaries including the one in our system, ensuring a consistent ability to assess the amount of task-relevant information in reconstructed images. Since adversarial training, which is grounded in the information bottleneck principle, demands the careful navigation of a dynamic trade-off between the transmission of task-relevant information and the mitigation of adversarial threats, this inevitable compromise inherently limits the amount of information that can be transmitted during standard communication, ultimately causing a pronounced decline in the receiver's accuracy. In stark contrast, our diffusion-based receiver demonstrates remarkable superiority, achieving significantly higher performance compared to both IBAL systems.

Fig. 5 shows how accuracy varies with SSIM and PSNR on the MNIST dataset. Here, "Direct" and "Permutation" mean the accuracy of $\bar{\Theta}$'s classification to the original images and permuted images that are directly noised by a certain level of Gaussian white noise [4] (as shown in Fig. 7), based on Eq. (12); we use a simple permutation method that divides the image into four parts and swaps the position of up-left corner and down-right corner parts. Essentially, our system's

---

[4]Here we just use those two baselines as comparisons to show that SSIM or PSNR may be disconnected to the semantic fidelity in some cases, and the appliance of permutation is beyond our scope of this paper. Future work may utilize permutation to further enhance the system performance against model inversion attack.
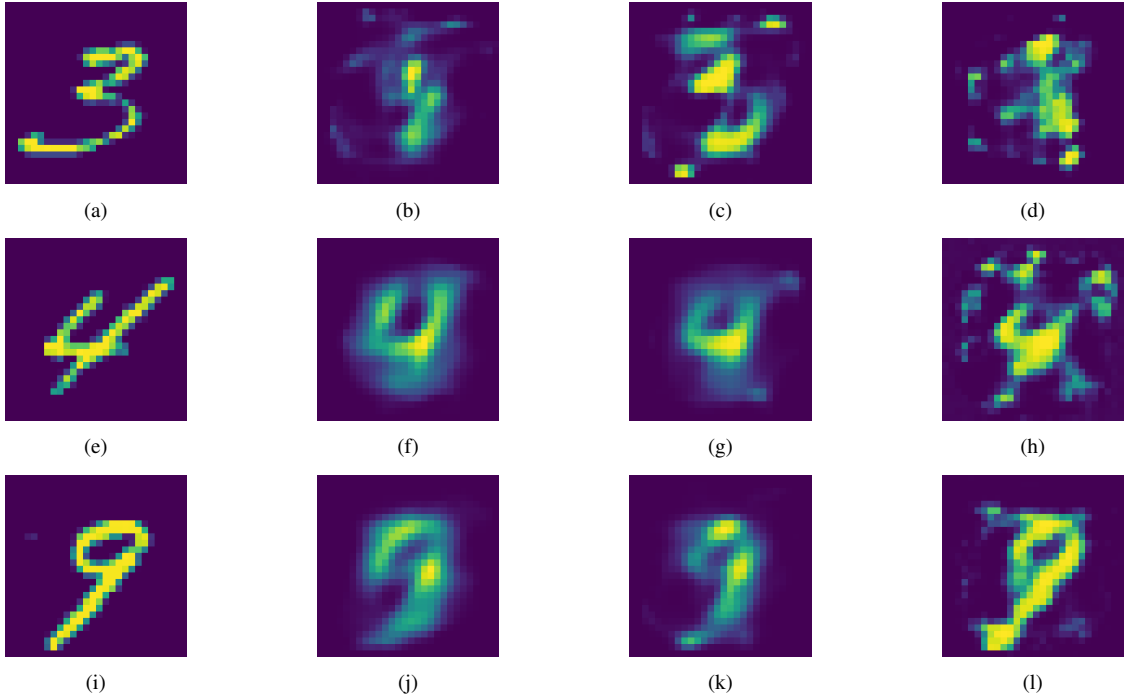
Fig. 6: Adversary's output images with different SNR or $t$ under MNIST dataset. 1st row: $-10$ dB; $t = 344$. 2nd row: 0 dB; $t = 184$. 3rd row: 10 dB; $t = 67$. (a)(e)(i) Original images. (b)(f)(j) IBAL system. (c)(g)(k) IBAL-D system. (d)(h)(l) Our system.
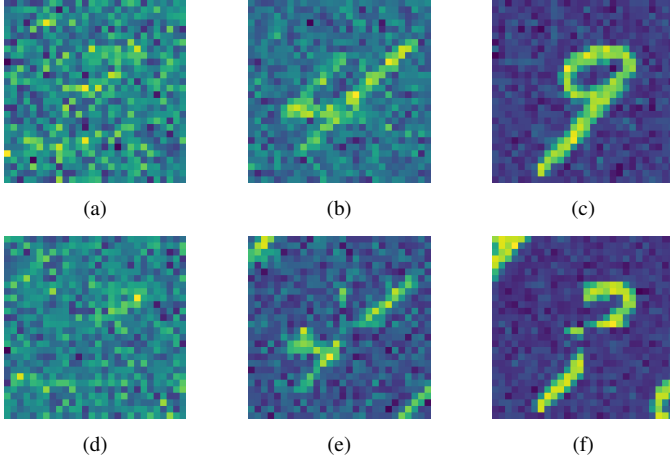


Fig. 7: (a)(b)(c) Direct mode and (d)(e)(f) permutation mode, with channel noise $-10$ dB, 0 dB, 10 dB. Original images are as (a)(e)(i) in Fig. 6.

is 0.78. This indicates that adversarial training merely distorts syntactic features without suppressing semantic information, since the key features remain intact or largely unchanged from the perspective of neural networks. It can be evidenced in Fig. 6 that the adversary's images in IBAL still have enough semantic information that can be used to categorize, even though their SSIM or PSNR values are lower. The lower plot in Fig. 5, showing the complicated correspondence between PSNR and accuracy, further supports the discussion above, since PSNR ignores perceptual quality and structural information that is important to a classification task [30]. In task-oriented semantic communication, the primary goal is to accurately reconstruct and understand the transmitted task-related information, rather than merely preserving pixel-level fidelity. Therefore, metrics like accuracy, which directly measures the system's ability to correctly interpret the semantic content, are more appropriate for evaluating the effectiveness of adversarial attacks.

The training strategy of our system prioritizes the recovery of semantic information. The adversary is only tasked with reconstructing images after the $\mathcal{T} - \mathcal{R}$ pair is fully trained to recover task-relevant semantic information, and the amount of information in $\mathcal{T}$'s output can be considered constant. Since we further process the adversary's output through a strong classifier that effectively converting the image reconstruction task into a classification task, it allows us to evaluate how much of the "task-relevant" information is contained within the reconstructed images, which is directly tied to the quality of the image reconstruction. However, due to the limited information provided by $\mathcal{T}$ and the fact that the $\mathcal{T}$'s output is not primarily designed for image reconstruction, the adversary's ability to recover syntactic information is
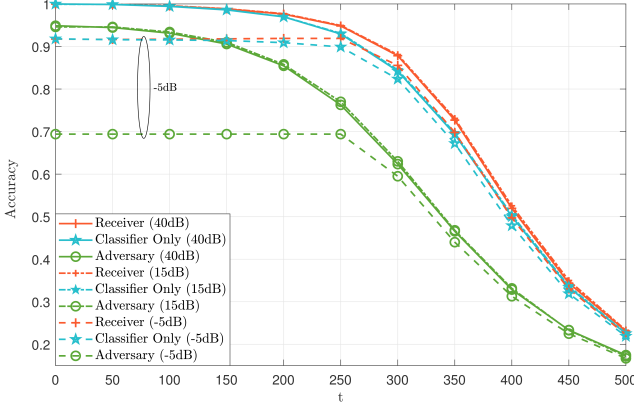
adversary $\mathcal{A}$, as well as the adversaries in the two IBAL systems, can partially achieve the recovery of the input $x$, as evidenced by their SSIM values being higher than those of the direct-noising scenario in most cases with the same accuracy. This aligns with the characteristic of adversaries that have denoising capabilities, as noted in [41]. Compared to the adversary in our system, despite the adversaries in both information-bottleneck-based systems effectively reduce the estimated SSIM and PSNR of the original input obtained by the attacker, we observe that there are cases that the adversary's classification accuracy can remain relatively high. For instance, the classification accuracy for our system's attacker (blue dashed line) at SSIM = 0.2 is 0.47, while for IBAL-D (purple dashed line with circle markers) at SSIM = 0.175, it

Fig. 8: System performances on MNIST dataset under AWGN channel with different channel SNRs.



Fig. 9: System performances using CIFAR10 dataset.

inherently biased, resulting in less "task-relevant" information embedded in the reconstructed images. If the $\mathcal{T} - \mathcal{A}$ pair is trained first, followed by the receiver $\mathcal{R}$, the training process would inevitably prioritize a "high-information task" (image reconstruction) before transitioning into a "low-information task" (classification). This approach would likely lead to the encoder producing excessive information, which, according to the evaluation standard defined in ***Def. (1)***, would benefit the adversary by providing more access to task-relevant semantic information, even if the recovered images have relatively low SSIM (or PSNR) values. Consequently, this strategy would compromise the privacy protection intended for the classifier.

Fig. 8 illustrates the robustness of our DiffSem system across diverse channel conditions. As outlined in Sec. III Part C, the transmitter dynamically adjusts the level of self-noising pre-added noise based on the estimated channel condition. When facing high SNR regimes, the self-noising module becomes the primary driver of signal distortion, meaning that variations in $t$ more directly influence the performance of both the receiver and the attacker, allowing precise control via $t$. Conversely, in poor channel conditions, the wireless channel itself dominates, allowing the self-noising module to flexibly control the quality of the transmitted signal while preventing over-compression to exert minimal impact on the signal state at the receiver. For instance, at a channel SNR of $-5$ dB, the dominance of channel noise results in a nearly constant information gap between the classifier and the adversary. We observe that for $t > 250$, the performance of both the receiver and the adversary slightly deteriorates compared to scenarios with higher channel SNR. This discrepancy may stem from the discrete nature of $\bar{\alpha}_t$, causing deviations between practical estimates and theoretical calculations. Overall, our self-noising module effectively modulates the degree of signal distortion in response to wireless channel conditions, maintaining an appropriate noise margin and ensuring reliable operation under both high and low SNR environments.

### C. Results on CIFAR10 dataset

We also implement the system on the CIFAR10 dataset, which has complex scenes that pose more challenges for semantic communication systems. Fig. 9 illustrates the re-
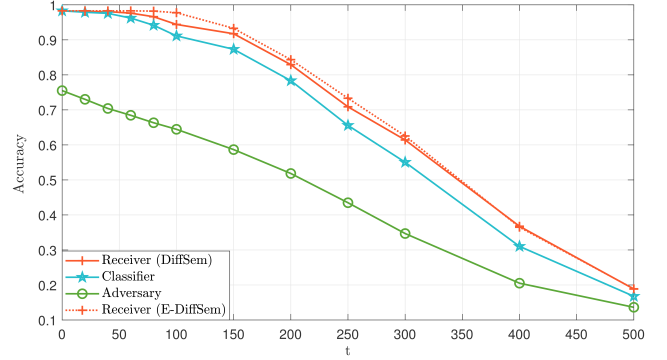
ceiver's classification accuracy and the attacker's success rate as functions of $t$. Given our assumption that the training of the $\mathcal{T} - \mathcal{R}$ pair takes precedence over the adversary's training, even at $t = 0$, the classification accuracy of the adversary's reconstructed images is not very high. This is because the $\mathcal{T} - \mathcal{R}$ pair's training prioritizes the completion of the label classification task, and the backpropagation process optimizes the $\mathcal{T}$'s neural network in a way that is not aligned with the adversary's image recovery requirements, which is acceptable since normal systems are not designed to meet attackers' needs.

The two red lines in Fig. 9 show that the system incorporating the diffusion module outperforms the classifier-only scenario. This indicates that imperfect labels can still provide coarse-grained constraints for the diffusion path. However, the additional performance improvement brought by E-DiffSem is less pronounced compared to the MNIST dataset. This discrepancy may be because the CIFAR10 dataset contains richer and more complex information than MNIST, which makes it more challenging for the imperfect labels estimated by the classifier to provide effective guidance to the diffusion module. In the future, we could explore integrating Transformer-based self-attention mechanisms or leveraging image semantic segmentation techniques to provide stronger guidance for the diffusion module, thereby further enhancing system performance.

As shown in Fig. 11, our system actively degrades the semantic utility of adversarial reconstructions. For instance, at $t = 150$ ($p = 2.3$ dB), the adversary's recovered images (Fig. 11(c)) retain only fragmented visual cues (e.g., partial textures), resulting in a classification accuracy of $58.7\%$, which is significantly below $87.3\%$ when bypassing diffusion, as well as $91.7\%$ and $93.3\%$ when using DiffSem and E-DiffSem. The gaps stem from our diffusion model's prioritized recovery of task-specific features (e.g., object categories) while discarding non-essential syntactic details. Such a design ensures that even if attackers partially reconstruct pixel-level content, they cannot infer semantic information effectively in a certain possibility.

Fig. 10 shows the accuracy versus SSIM and PSNR on the CIFAR10 dataset. It highlights the weak correlation between traditional image quality metrics and semantic security. For instance, when SSIM is $0.12$, the attacker's category accuracy still reaches $51.8\%$, whereas in the direct mode, with SSIM
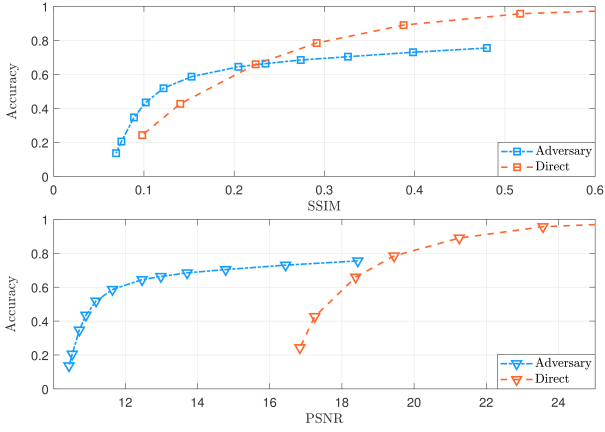
Fig. 10: Plot of accuracy versus SSIM and PSNR using CIFAR10 dataset.

equal to $0.14$, the accuracy is only $42.6\%$. This difference arises from the attacker's ability to implicitly learn class-discriminative features, while variations in pixel texture are insufficient to prevent the attacker from acquiring these features. Other curves in this figure can further demonstrate that SSIM and PSNR are not suitable metrics for evaluating the amount of task-relevant information obtained by the attacker in task-oriented communication systems.

As a supplement to the discussion, here we remark the computational complexity of our proposed system. Our hardware conditions are: CPU, Intel Xeon Gold 5220; GPU, GeForce RTX 3090. Although the diffusion model introduces additional computational overhead, the DDIM acceleration strategy restricts the additional single-sample inference time of DiffSem to approximately $0.1$s per batch (batch size 256) on CIFAR10. It shows that the performance enhancement does not impose a substantial computational burden given the relatively small dataset. For higher-resolution datasets, such as ImageNet, the trade-off between performance improvement and computational cost may become a critical consideration, especially in real-time applications where latency is a key factor. Future work could focus on optimizing the diffusion process, potentially exploring techniques such as knowledge distillation, model pruning, or adopting more efficient neural network architectures to address this challenge.

## V. CONCLUSION

In this paper, we have introduced a diffusion-based semantic communication framework named DiffSem, to address the dual challenges of optimizing task performance and ensuring privacy protection in semantic communication systems. By leveraging the diffusion mechanism and integrating label embedding, DiffSem effectively reconstructs semantic information and significantly improves task accuracy, as evidenced by a $10.3\%$ boost on the MNIST dataset. To tackle the growing concern of privacy risks posed by model inversion attacks, we proposed a novel evaluation metric that quantifies the semantic fidelity of adversarial outputs, offering a precise, task-relevant method to assess privacy vulnerabilities. Furthermore, we developed a non-adversarial privacy protection strategy that prioritizes the preservation of semantic fidelity for legitimate

receivers while constraining the semantic information accessible to adversaries. This approach, combined with an adaptive noise injection mechanism, enables the system to maintain a clear distinction between semantic and syntactic information recovery, thereby advancing the understanding of their respective roles in privacy and communication. Experimental results also highlight the limitations of traditional image quality metrics, such as PSNR and SSIM, in reflecting task-relevant semantic fidelity, highlighting the utility of our proposed metric. Overall, this work provides a unified framework that balances task optimization and privacy protection, contributing to the development of secure and efficient semantic communication systems.

## APPENDIX

Here we derive the Eq. (25) in detail. Since we have

$$\mathbb{E}\left(\|z'\|^2\right) = \bar{\alpha}_{T'}\|f_0\|^2 + (1 - \bar{\alpha}_{T'}) D \tag{29}$$

and

$$\begin{aligned} \mathbb{E}\left(\|\hat{z}''\|^2\right) &= \bar{\alpha}_{T'}\|f_0\|^2 + (1 - \bar{\alpha}_{T''}) D + \sigma_{\text{ch}}^2 D \\ &= \bar{\alpha}_{T'}\|f_0\|^2 + (1 - \bar{\alpha}_{T''}) D \\ &\quad + \frac{\bar{\alpha}_{T'}\mathbb{E}\left(\|f_0\|^2\right) + (1 - \bar{\alpha}_{T''})D}{\gamma} D, \end{aligned} \tag{30}$$

we can get

$$1 - \bar{\alpha}_{T''} + \frac{\bar{\alpha}_{T'}\mathbb{E}\left(\|f_0\|^2\right) + (1 - \bar{\alpha}_{T''})D}{\gamma} = 1 - \bar{\alpha}_{T'} \tag{31}$$

when $\mathbb{E}\left(\|z'\|^2\right) = \mathbb{E}\left(\|\hat{z}''\|^2\right)$. Then

$$\bar{\alpha}_{T''}\left(1 + \frac{D}{\gamma}\right) = \bar{\alpha}_{T'}\left(1 + \frac{\mathbb{E}(\|f_0\|^2)}{\gamma}\right) + \frac{D}{\gamma}. \tag{32}$$

Arranging the formula above, we can get

$$\bar{\alpha}_{T''} = \frac{\bar{\alpha}_{T'}\left(\gamma + \mathbb{E}\left(\|f_0\|^2\right)\right) + D}{\gamma + D}. \tag{33}$$

## REFERENCES

[1] W. Weaver, "Recent contributions to the mathematical theory of communication," *ETC: a review of general semantics*, pp. 261–281, 1953.
[2] Y. Shao, Q. Cao, and D. Gündüz, "A theory of semantic communication," *IEEE Transactions on Mobile Computing*, 2024.
[3] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE transactions on signal processing*, vol. 69, pp. 2663–2675, 2021.
[4] E. Beck, C. Bockelmann, and A. Dekorsy, "Semantic Communication: An Information Bottleneck View," *ArXiv*, vol. abs/2204.13366, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID: 248427116
[5] L. Yan, Z. Qin, R. Zhang, Y. Li, and G. Y. Li, "Resource allocation for text semantic communications," *IEEE Wireless Communications Letters*, vol. 11, no. 7, pp. 1394–1398, 2022.
[6] X. Peng, Z. Qin, X. Tao, J. Lu, and L. Hanzo, "A robust semantic text communication system," *IEEE Transactions on Wireless Communications*, 2024.
[7] Z. Lyu, G. Zhu, J. Xu, B. Ai, and S. Cui, "Semantic communications for image recovery and classification via deep joint source and channel coding," *IEEE Transactions on Wireless Communications*, vol. 23, no. 8, pp. 8388–8404, 2024.
[8] H. Wu, Y. Shao, C. Bian, K. Mikolajczyk, and D. Gündüz, "Deep joint source-channel coding for adaptive image transmission over MIMO channels," *IEEE Transactions on Wireless Communications*, 2024.
[9] Q. Ma, J. Pan, and C. Bai, "Direction-oriented visual-semantic embedding model for remote sensing image-text retrieval," *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
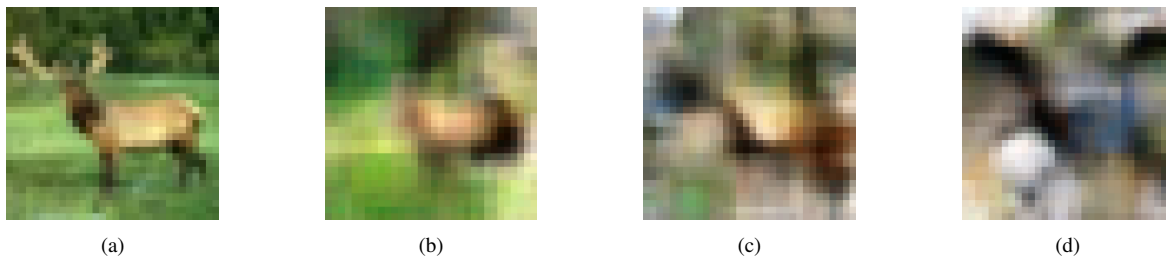
Fig. 11: Adversary's output images with different $t$ under CIFAR10 dataset. (a) Original image. (b)(c)(d) Adversary's recovery when $t = 60$, 150 and 300.

[10] H. Li, H. Tong, S. Wang, N. Yang, Z. Yang, and C. Yin, "Video semantic communication with major object extraction and contextual video encoding," in *2024 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2024, pp. 1–6.

[11] S. Wang, J. Dai, Z. Liang, K. Niu, Z. Si, C. Dong, X. Qin, and P. Zhang, "Wireless deep video semantic transmission," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 214–229, 2022.

[12] Y. Cang, M. Chen, Z. Yang, Y. Hu, Y. Wang, C. Huang, and Z. Zhang, "Online resource allocation for semantic-aware edge computing systems," *IEEE Internet of Things Journal*, vol. 11, no. 17, pp. 28 094–28 110, 2023.

[13] W. Yang, Z. Q. Liew, W. Y. B. Lim, Z. Xiong, D. Niyato, X. Chi, X. Cao, and K. B. Letaief, "Semantic communication meets edge intelligence," *IEEE wireless communications*, vol. 29, no. 5, pp. 28–35, 2022.

[14] K. Muhammad, T. Hussain, H. Ullah, J. Del Ser, M. Rezaei, N. Kumar, M. Hijji, P. Bellavista, and V. H. C. de Albuquerque, "Vision-based semantic segmentation in scene understanding for autonomous driving: Recent achievements, challenges, and outlooks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 22 694–22 715, 2022.

[15] G. Zheng, Q. Ni, K. Navaie, H. Pervaiz, and C. Zarakovitis, "A distributed learning architecture for semantic communication in autonomous driving networks for task offloading," *IEEE Communications Magazine*, vol. 61, no. 11, pp. 64–68, 2023.

[16] H. Zhang, H. Wang, Y. Li, K. Long, and V. C. Leung, "Toward intelligent resource allocation on task-oriented semantic communication," *IEEE Wireless Communications*, vol. 30, no. 3, pp. 70–77, 2023.

[17] D. Gündüz, Z. Qin, I. E. Aguerri, H. S. Dhillon, Z. Yang, A. Yener, K. K. Wong, and C.-B. Chae, "Beyond transmitting bits: Context, semantics, and task-oriented communications," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 5–41, 2022.

[18] M. Jankowski, D. Gündüz, and K. Mikolajczyk, "Deep joint source-channel coding for wireless image retrieval," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 5070–5074.

[19] H. Sun, W. Ni, H. Tian, J. Zheng, G. Nie, and D. Niyato, "A Hybrid Federated Learning Framework for Task-Oriented Semantic Communication," *IEEE Internet of Things Journal*, 2025.

[20] J. Peng, H. Xing, Y. Li, L. Feng, L. Xu, and X. Lei, "Task-oriented multi-user semantic communication with lightweight semantic encoder and fast training for resource-constrained terminal devices," *IEEE Wireless Communications Letters*, 2024.

[21] C. Cai, X. Yuan, and Y.-J. A. Zhang, "End-to-End Learning for Task-Oriented Semantic Communications Over MIMO Channels: An Information-Theoretic Framework," *IEEE Journal on Selected Areas in Communications*, 2025.

[22] Z. Yang, M. Chen, G. Li, Y. Yang, and Z. Zhang, "Secure semantic communications: Fundamentals and challenges," *IEEE network*, 2024.

[23] N.-B. Nguyen, K. Chandrasegaran, M. Abdollahzadeh, and N.-M. Cheung, "Re-thinking model inversion attacks against deep neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16 384–16 393.

[24] T. Ha, T. K. Dang, T. T. Dang, T. A. Truong, and M. T. Nguyen, "Differential privacy in deep learning: an overview," in *2019 International Conference on Advanced Computing and Applications (ACOMP)*. IEEE, 2019, pp. 97–102.

[25] K. Pan, Y.-S. Ong, M. Gong, H. Li, A. K. Qin, and Y. Gao, "Differential privacy in deep learning: A literature survey," *Neurocomputing*, p. 127663, 2024.

[26] L. Zhao, D. Wu, and L. Zhou, "Data utilization versus privacy protection in semantic communications," *IEEE Wireless Communications*, vol. 30, no. 3, pp. 44–50, 2023.

[27] T. Xiao, Y.-H. Tsai, K. Sohn, M. Chandraker, and M.-H. Yang, "Adversarial learning of privacy-preserving and task-oriented representations," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 434–12 441.

[28] A. Li, Y. Duan, H. Yang, Y. Chen, and J. Yang, "TIPRDC: task-independent privacy-respecting data crowdsourcing framework for deep learning with anonymized intermediate representations," in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 824–832.

[29] Y. Wang, S. Guo, Y. Deng, H. Zhang, and Y. Fang, "Privacy-preserving task-oriented semantic communications against model inversion attacks," *IEEE Transactions on Wireless Communications*, 2024.

[30] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[31] G. Shi, D. Gao, X. Song, J. Chai, M. Yang, X. Xie, L. Li, and X. Li, "A new communication paradigm: From bit accuracy to semantic fidelity," *arXiv preprint arXiv:2101.12649*, 2021.

[32] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 9, pp. 10 850–10 869, 2023.

[33] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, and M.-H. Yang, "Diffusion models: A comprehensive survey of methods and applications," *ACM Computing Surveys*, vol. 56, no. 4, pp. 1–39, 2023.

[34] E. Grassucci, C. Marinoni, A. Rodriguez, and D. Comminiello, "Diffusion models for audio semantic communication," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 13 136–13 140.

[35] T. Wu, Z. Chen, D. He, L. Qian, Y. Xu, M. Tao, and W. Zhang, "CDDM: Channel denoising diffusion models for wireless semantic communications," *IEEE Transactions on Wireless Communications*, 2024.

[36] L. Guo, W. Chen, Y. Sun, B. Ai, N. Pappas, and T. Quek, "Diffusion-driven semantic communication for generative models with bandwidth constraints," *IEEE Transactions on Wireless Communications*, 2025.

[37] E. Grassucci, S. Barbarossa, and D. Comminiello, "Generative semantic communication: Diffusion models beyond bit recovery," *arXiv preprint arXiv:2306.04321*, 2023.

[38] H. Fang, Y. Qiu, H. Yu, W. Yu, J. Kong, B. Chong, B. Chen, X. Wang, S.-T. Xia, and K. Xu, "Privacy leakage on dnns: A survey of model inversion attacks and defenses," *arXiv preprint arXiv:2402.04013*, 2024.

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

[40] Y. Chen, Q. Yang, Z. Shi, and J. Chen, "The model inversion eavesdropping attack in semantic communication systems," in *GLOBECOM 2023-2023 IEEE Global Communications Conference*. IEEE, 2023, pp. 5171–5177.

[41] X. Wang and L. Lu, "Implementation of DNN-based physical-layer network coding," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 6, pp. 7380–7394, 2023.