

# The Model Inversion Eavesdropping Attack in Semantic Communication Systems

Yuhao Chen, Qianqian Yang<sup>†</sup>, Zhiguo Shi, Jiming Chen

College of Control Science and Engineering, Zhejiang University, Hangzhou 310007, China

College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310007, China

The State Key Laboratory of Industrial Control Technology, Hangzhou 310007, China

{csechenyh, qianqianyang20<sup>†</sup>, shizg, cjm}@zju.edu.cn

**Abstract**—In recent years, semantic communication has been a popular research topic for its superiority in communication efficiency. As semantic communication relies on deep learning to extract meaning from raw messages, it is vulnerable to attacks targeting deep learning models. In this paper, we introduce the model inversion eavesdropping attack (MIEA) to reveal the risk of privacy leaks in the semantic communication system. In MIEA, the attacker first eavesdrops the signal being transmitted by the semantic communication system and then performs model inversion attack to reconstruct the raw message, where both the white-box and black-box settings are considered. Evaluation results show that MIEA can successfully reconstruct the raw message with good quality under different channel conditions. We then propose a defense method based on random permutation and substitution to defend against MIEA in order to achieve secure semantic communication. Our experimental results demonstrate the effectiveness of the proposed defense method in preventing MIEA.

## I. Introduction

Recently, semantic communication has been widely believed to be one of the core technologies for the sixth generation (6G) of wireless networks because of its high communication efficiency [?]. Compared with the current research on communication which focuses on transmitting mapped bit sequences of the raw message [?], [?], [?], semantic communication systems transmit compacted semantic features. Existing literature in semantic communication mainly exploits the deep learning (DL) techniques to extract the semantic features from the raw message. For instance, Han et al. [?] proposed to extract the text-related features from the speech signal as the semantic features and remove the redundant content. On the receiver's side, the semantic features can be reconstructed by a deep learning model into the original message or directly applied for downstream tasks such as image classification and speech recognition.

Although many works have been proposed for semantic communication considering different aspects, few studies have taken into account the security problems [?], [?], [?]. Tung et al. [?] proposed to encrypt the transmitted signal in semantic communication, but the encryption algorithm incurs a large computation overhead. Security is crucial in semantic communication for two main reasons. Firstly, semantic communication is more prone to privacy leakage compared to traditional communication. In traditional

communication systems, the bit sequences being transmitted contain redundant bits to ensure reliable transmission, which can be used to provide a certain level of privacy protection. However, the semantic communication systems transmit compact and more semantic-related symbols which may reveal more private information. Secondly, deep-learning-based semantic communication may be vulnerable to attacks targeting DL models. Extensive studies have been conducted on attacks on the DL model, a review of which can be referred to [?]. If the semantic features being transmitted are eavesdropped by a malicious attacker, the attacker can reconstruct the raw message by utilizing the DL-based attack techniques. The attacker can also add perturbation to the transmitted data, causing the semantic communication system to make incorrect decisions on downstream tasks. For example, Sagduyu et al. [?] proposed a multi-domain evasion attack to cause the semantic communication system to make incorrect classifications, which is achieved by introducing noises to input images or the semantic features. Du et al. [?] proposed a semantic data poisoning attack, which causes the receiver to receive irrelevant messages from the transmitter. For example, the receiver wants to receive an image with a pear but gets an image with an apple instead. This attack is performed by minimizing the difference between the semantic features of the targeted message and the irrelevant message.

In this paper, we consider the security issue in semantic communication systems and introduce the model inversion eavesdropping attack (MIEA) for semantic communication, where an attacker eavesdrops the transmitted symbols and attempts to reconstruct the original message from them by inverting the DL model used at the transmitter. We perform MIEA under both the white-box and the black-box settings. The attacker has knowledge of the DL model in the white-box setting while not in the black-box setting. To defend against MIEA, we also propose a defense method based on random permutation and substitution. Evaluations demonstrate that the MIEA attack works under different channel conditions, i.e., different values of the signal-to-noise ratio (SNR), which reveals the risk of privacy leaks in semantic transmission. Numerical results also validate the effectiveness of our proposed

defense method.

This paper is organized as follows: Section ?? introduces the basic ideas of semantic communications. In section ??, we present the proposed MIEA under both the white-box and black-box setting, and propose our defense method. In section ??, we evaluate the effectiveness of the proposed MIEA and the proposed defense method. Section ?? concludes our work.

## II. Fundamentals

In this section, we provide the fundamentals of semantic communication and the eavesdropping performed by the attacker. We consider a semantic communication system which transmits images over wireless channels. As shown in Fig. ??, the transmitter of the semantic communication system consists of a semantic encoder and a channel encoder. The semantic encoder extracts the semantic features  $\mathbf{z}$  from the raw image  $\mathbf{x}$ , while the channel encoder maps  $\mathbf{z}$  into the transmitted features  $\mathbf{y}_f \in \mathbb{R}^{h \times w \times c}$ , where  $h, w, c$  denote the height, the width and the channel of the transmitted features respectively. Before transmission,  $\mathbf{y}_f$  is reshaped into the transmitted symbols  $\mathbf{y} \in \mathbb{R}^{N \times 2}$ , where  $N = \frac{h \times w \times c}{2}$  and the two channels are the real parts and imaginary parts of the signal to be transmitted, respectively.  $\mathbf{y}$  is then transmitted over a wireless channel, which we denote as the main channel to distinguish from the channel used by the attacker. The received signal  $\hat{\mathbf{y}}$  at the receiver side can be characterized by

$$\hat{\mathbf{y}} = \mathbf{H}_m \mathbf{y} + \mathbf{n}_m, \quad (1)$$

where  $\mathbf{H}_m$  is a matrix which reflects the main channel effect such as multi-path propagation, fading and interference, while  $\mathbf{n}_m$  is a zero-mean additive white Gaussian noise. The receiver of the semantic communication system consists of a channel decoder and a semantic decoder. The receiver first reshapes  $\hat{\mathbf{y}}$  back to the transmitted features  $\hat{\mathbf{y}}_f$ . Then the channel decoder maps  $\hat{\mathbf{y}}_f$  back to the semantic features  $\hat{\mathbf{z}}$ . The semantic decoder then reconstructs the image  $\hat{\mathbf{x}}$  from  $\hat{\mathbf{z}}$ . We jointly train the semantic encoder, channel encoder, semantic decoder and channel decoder using the following loss function:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \lambda T(\hat{\mathbf{x}}), \quad (2)$$

where  $N$  is the number of the training data batch and

$$T(\hat{\mathbf{x}}) = \sum_{i,j} (|\hat{\mathbf{x}}_{i+1,j} - \hat{\mathbf{x}}_{i,j}|^2 + |\hat{\mathbf{x}}_{i,j+1} - \hat{\mathbf{x}}_{i,j}|^2)^{\beta/2}. \quad (3)$$

The first term in (??) computes the mean square error (MSE) between  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ . The second term  $T(\hat{\mathbf{x}})$  is the total variation [?] that measures the smoothness of the reconstructed image  $\hat{\mathbf{x}}$ , where  $\hat{\mathbf{x}}_{i,j}$  denotes the pixel value at the position  $(i, j)$  and  $\beta$  controls the smoothness of the image, with larger  $\beta$  being more piecewise-smooth. The

hyper-parameter  $\lambda$  balances the two terms. In our work, we choose  $\beta = 1$  and  $\lambda = 1$ .

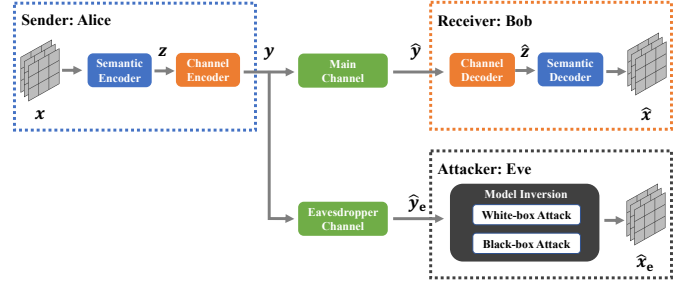


Fig. 1. Illustration of the semantic communication with MIEA.

Next, we introduce how an attacker eavesdrops the transmitted signal under the semantic communication system. We follow the naming convention in the security research, with Alice, Bob and Eve representing the sender, receiver and attacker respectively. Suppose Alice wants to send an image to Bob. As shown in the lower part of Fig. ??, since the transmitted symbols  $\mathbf{y}$  is transmitted over the wireless channel, it can easily be captured by any unauthorized receiver. Assume that there exists an attacker Eve who intercepts  $\mathbf{y}$  and attempts to reconstruct the raw image from it. The wireless channel between Alice and Eve is referred to as the eavesdropper channel [?]. The received signal at Eve is given by

$$\hat{\mathbf{y}}_e = \mathbf{H}_e \mathbf{y} + \mathbf{n}_e, \quad (4)$$

Similarly,  $\mathbf{H}_e(\cdot)$  represents the eavesdropper channel matrix and  $\mathbf{n}_e$  is a zero-mean additive white Gaussian noise. After eavesdropping  $\hat{\mathbf{y}}_e$ , Eve is able to reconstruct the image, denoted as  $\hat{\mathbf{x}}_e$ , which will be detailed section ??. Note that to avoid confusion, we use the received image to denote  $\hat{\mathbf{x}}$  received by Bob and the eavesdropped image to denote  $\hat{\mathbf{x}}_e$  eavesdropped by Eve.

## III. The Proposed MIEA and its Defense

In this section, we first elaborate the idea of MIEA. To reconstruct  $\hat{\mathbf{x}}_e$ , Eve performs MIA [?] using either the white-box attack or the black-box attack, which depends on the knowledge of the semantic encoder and channel encoder that Eve has. Then we propose an effective defense method that defends against both types of attack.

### A. White-box Attack

In the white-box attack, Eve knows the parameters and structure of the semantic encoder and channel encoder. For example, the semantic communication system is publicly available or available through purchase, such as JPEG. In this case, Eve can directly use the semantic encoder and channel encoder to reconstruct the image. We denote the two encoders as a single function  $f(\cdot)$  which maps a given image  $\mathbf{x}$  to the transmitted symbol  $\mathbf{y}$ , that is,  $\mathbf{y} = f(\mathbf{x})$ .

The reconstructed image  $\hat{\mathbf{x}}_e$  can be obtained by solving the following optimization problem:

$$\hat{\mathbf{x}}_e = \arg \min_{\mathbf{x}} \|\hat{\mathbf{y}}_e - f(\mathbf{x})\|^2 + \lambda T(\mathbf{x}). \quad (5)$$

The first term in (5) is the MSE between  $\hat{\mathbf{y}}_e$  and  $f(\mathbf{x})$ , while the second term  $T(\mathbf{x})$  is the total variation defined in (2) that guarantees the smoothness of  $\mathbf{x}$ . Similar to (2), we set  $\beta = 1$  and  $\lambda = 1$  here. The optimization problem (5) can be solved by performing the gradient descent, which iteratively updates the input  $\mathbf{x}$  so that (5) can be minimized.

### B. Black-box Attack

In the black-box attack, Eve lacks knowledge of the parameters and structures of both encoders. In this case, Eve uses an inverse network of the two encoders, denoted as  $f^{-1}(\cdot)$ , to inverse  $\hat{\mathbf{y}}_e$  back to  $\hat{\mathbf{x}}_e$ . Specifically,  $f^{-1}$  takes  $\hat{\mathbf{y}}_e$  as input and outputs  $\mathbf{x}$ , i.e.,  $f^{-1}(\hat{\mathbf{y}}_e) = \mathbf{x}$ . To train  $f^{-1}$ , we assume that Eve can feed a batch of samples  $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$  into the encoder and capture the corresponding transmitted symbols  $\mathbb{Y}_e = \{\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_m\}$ , where  $m$  is the number of the samples. Eve then trains  $f^{-1}$  using  $\mathbb{Y}_e$  as the input and  $\mathbb{X}$  as the ground truth output. We use the  $l_2$  norm as the loss function and employ stochastic gradient descent to train the inverse network:

$$f^{-1} = \arg \min_g \frac{1}{m} \sum_{i=1}^m \|g(\hat{\mathbf{y}}_i) - \mathbf{x}_i\|^2, \quad (6)$$

where we use  $g$  to represent the inverse network being optimized. Once the inverse network is trained, Eve is able to reconstruct the image from any newly eavesdropped signal.

### C. Defense Method against MIEA

To defend against MIA in deep learning, researchers have proposed various techniques such as differential privacy (DP) [?] or attacker-aware training [?]. The DP technique involves adding an additive Laplacian noise to the raw image, while the attacker-aware training adds a regularization term to the loss function during training, which maximizes the MSE between the reconstructed image and the raw image. Both methods prevent Eve from reconstructing high-quality images while maintaining the performance of downstream tasks. However, in the transmission task considered in this paper, both Bob and Eve attempt to reconstruct the image from  $\hat{\mathbf{y}}$  and  $\hat{\mathbf{y}}_e$  respectively. Furthermore, we assume that both the main channel and eavesdropper channel are AWGN channel, i.e.,  $H_m = H_e$ . Then the difference between  $\hat{\mathbf{y}}$  and  $\hat{\mathbf{y}}_e$  is  $|\mathbf{n}_m - \mathbf{n}_e|$ , which is relatively small. Therefore, if Eve fails to reconstruct high-quality images under the defense methods above, Bob will also fail, which contradicts the goal of the transmission task. To prevent Eve from reconstructing images as well as maintaining the image quality received by Bob, an intuitive solution is to encrypt  $\mathbf{y}$  using

common cryptography algorithms, but this will incur a large computation overhead. To reduce the computation overhead, we propose a defense method based on random permutation and substitution of the transmitted features  $\hat{\mathbf{y}}_f$ , which can simultaneously defend against both types of attack.

**Random Permutation and Substitution.** We first introduce the random permutation operation. For the transmitted features  $\mathbf{y}_f$ , we randomly permute the tensor along the first dimension  $h$ . We define the permutation scheme  $P$  as a random permutation of the array  $[0, 1, \dots, h-1]$ , where each element represents the index (0-indexed) of  $\mathbf{y}_f$ . After applying  $P$ , we obtain the permuted transmitted features  $\mathbf{y}_f^P$ .

Then we perform the substitution operation on  $\mathbf{y}_f^P$  by swapping some of the  $y_i$  with  $y'_i$  from another transmitted features  $\mathbf{y}'_f$ , where  $i$  in  $y_i$  and  $y'_i$  indicates that substitution is performed in the same position of both transmitted features. Note that we remove the subscript  $f$  in  $y_i$  to avoid complex notations. If Alice sends several images to Bob,  $\mathbf{y}'_f$  can be the transmitted features of the next image. If Alice only sends one image to Bob,  $\mathbf{y}'_f$  can be a random-noise tensor, which will also be sent to Bob. Similarly, we define the substitution scheme  $S$  as a sub-array of the array  $[0, 1, \dots, h-1]$ , where every  $i$  in  $S$  indicates the  $y_{f,i}$  that should be substituted. After substitution,  $\mathbf{y}_f^P$  becomes  $\mathbf{y}_f^S$ .

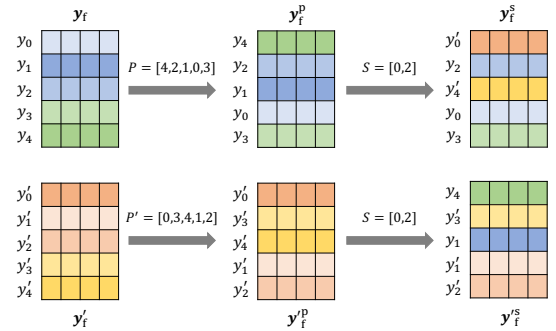


Fig. 2. An example of the proposed defense method.

We give an example to explain the idea of random permutation and substitution. As shown in Fig. ??, assume that  $h = 5$ , then  $\mathbf{y}_f = [y_0, y_1, \dots, y_4]$ , where  $y_i \in \mathbb{R}^{1 \times w \times c}$ ,  $i = 0, 1, \dots, 4$ . We also assume that there is another transmitted features  $\mathbf{y}'_f = [y'_0, y'_1, \dots, y'_4]$ . If  $P = [4, 2, 1, 0, 3]$  for  $\mathbf{y}_f$  and  $P' = [0, 3, 4, 1, 2]$  for  $\mathbf{y}'_f$ , then  $\mathbf{y}_f^P = [y_4, y_2, y_1, y_0, y_3]$  and  $\mathbf{y}_f'^P = [y'_0, y'_3, y'_4, y'_1, y'_2]$ . If  $\mathbf{y}_f$  and  $\mathbf{y}'_f$  share a common  $S = [0, 2]$ , then  $\mathbf{y}_f^S = [y'_0, y_2, y'_4, y_0, y_3]$  and  $\mathbf{y}_f'^S = [y_4, y'_3, y_1, y'_1, y'_2]$ . Suppose that Bob knows the  $P$  and  $S$  before transmission. After reshaping  $\hat{\mathbf{y}}$ , Bob will first recover  $\hat{\mathbf{y}}_f$  from  $\hat{\mathbf{y}}_f^S$  and then feed  $\hat{\mathbf{y}}_f$  into the channel decoder. However, since Eve does not know  $P$  and  $S$ , Eve will try to reconstruct  $\hat{\mathbf{x}}_e$  directly from  $\hat{\mathbf{y}}_f^S$ , which is shown to be infeasible in ?? . Moreover, since different  $P$  and  $S$  are used for each transmission, it would be difficult for Eve

to determine the correct  $P$  and  $S$  for each eavesdropped signal.

**Scheme Selection.** The proposed method is dependent on Bob having knowledge of  $P$  and  $S$  before transmission. Hence it is necessary for Alice and Bob to share two common sets of schemes, namely the permutation scheme set  $\mathbb{P}$  and the substitution set  $\mathbb{S}$ , which are kept secret from Eve. Both sets comprise multiple schemes that can be employed for permutation and substitution. Before each image transmission, Alice generates a value pair  $V = \{p, s\}$ , which is used to select the corresponding  $P$  and  $S$  from  $\mathbb{P}$  and  $\mathbb{S}$ , respectively.  $V$  is first encrypted using a secret key  $K$  shared between Alice and Bob. Then the encrypted  $V$  is transmitted to Bob, which cannot be modified by the main channel. Hence error-free techniques such as error correction and retransmission are utilized to transmit  $V$ . After receiving the  $\hat{y}_s$  and the encrypted  $V$ , Bob decrypts  $V$  using  $K$  and determines  $P$  and  $S$ , from which  $\hat{y}$  can be recovered.

#### IV. Evaluations

In this section, we present our experiments to evaluate MIEA and the proposed defense method. We first evaluate MIEA's performance for the white-box attack and black-box attack. We then show the effectiveness of the proposed defense method. We use the semantic communication model DeepJSCC [?] to transmit images from the CelebA dataset [?], which we crop and resize to  $180 \times 180$  in the evaluation. The semantic encoder and decoder each have four convolutional layers, while the channel encoder and decoder have one convolutional layer. We assume both the main channel and eavesdropper channel to be AWGN channel and denote the channel condition as the combination of the main channel's SNR and the eavesdropper channel's SNR. Although the main channel is not considered in MIEA, we still perform evaluation under different SNRs of the main channel, as we use different DeepJSCC models for each SNR value. For each evaluation, we consider the SNR of both channels to be 0dB, 10dB and 20dB, resulting in nine different channel conditions.

##### A. Evaluation Setup

Before evaluating the performance of both attacks, we train the DeepJSCC model on the CelebA dataset using three different SNR values for the main channel (0dB, 10dB, and 20dB), resulting in three distinct DeepJSCC models. As stated in [?], the SNR value determines the standard deviation of  $\mathbf{n}_m$  when the transmission power is normalized to 1. We train the CelebA dataset with a batch size of 128, using Adam [?] as the optimizer with a learning rate of  $10^{-3}$ .

To measure the image quality, we use two metrics, i.e., the structural similarity index measure (SSIM) and the peak signal-to-noise ratio (PSNR) [?], where higher values of SSIM and PSNR indicate better quality.

##### B. Evaluation of MIEA

We first evaluate MIEA for the two types of attack. For the white-box attack, where Eve reconstructs the image by minimizing (??), we employ Adam [?] as the optimizer with a learning rate of  $10^{-3}$  and we initialize  $\hat{x}_e$  to an all-zero tensor. For the black-box attack, we use an inverse network  $f^{-1}(\cdot)$  consisting of an upsampling layer and two convolution layers. Then we train  $f^{-1}(\cdot)$  by solving the optimization problem in (??), where we choose the CelebA test dataset as  $\mathbb{X}$  and obtain its corresponding transmitted symbols  $\mathbb{Y}$ . Similarly, we use Adam as the optimizer and set the learning rate to  $10^{-3}$ .

Fig. ?? shows the performance of MIEA on the DeepJSCC model under different channel conditions, with the SSIM and PSNR given below each image. The first two columns in Fig. ?? are baselines for comparisons with the images eavesdropped by MIEA, where the first column displays the original images  $\mathbf{x}$  transmitted by Alice, and the second column shows the images received by Bob under different main channel's SNRs. As shown in the first two columns, increasing the SNR improves image quality, as indicated by higher average SSIM and PSNR values. Additionally, higher SNRs reveal more details in the images, such as the female's hair.

The remaining columns in Fig. ?? display the eavesdropped images obtained by MIEA. In the following evaluations in this paper, for each channel condition, we show the eavesdropped image by the white-box attack on the left and the one obtained by the black-box attack on the right. Additionally, Table?? lists the average SSIM and PSNR of the eavesdropped images of individuals selected from the CelebA training set. The first row in the table shows the quality of the images received by Bob, and for each channel condition, the values on the top show the quality of eavesdropped images obtained by the white-box attack, and the values on the bottom show the quality of the eavesdropped images obtained by the black-box attack. We note that we choose CelebA training set for evaluation because the CelebA test set is used for training  $f^{-1}(\cdot)$  in the black-box attack. It can be seen from Fig. ?? and Table ?? that the quality of eavesdropped images improves as the SNR of the eavesdropper channel increases for a given SNR of the main channel. Moreover, for a given SNR of the eavesdropper channel, the quality of eavesdropped images is similar under different SNRs of the main channel. It also can be observed that the SSIM and PSNR values in the black-box attack are generally larger than those in the white-box attack. This is because the black-box attack requires training  $f^{-1}(\cdot)$  before reconstructing any image from the eavesdropped signal, which needs many samples from  $\mathbb{X}$  and  $\mathbb{Y}$ . In contrast, the white-box attack directly reconstructs the image from the eavesdropped signal without any training in advance. Although the SSIM and the PSNR of the eavesdropped images in both attacks are lower than those

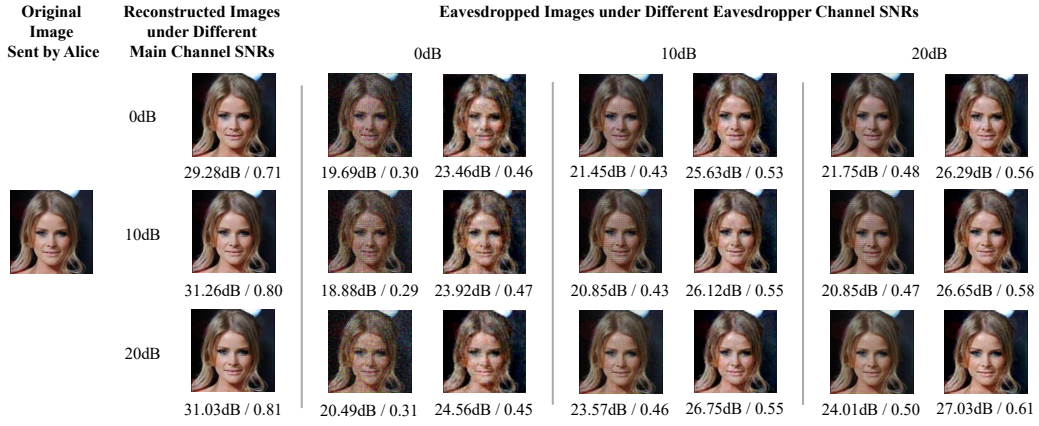


Fig. 3. Visualization of MIEA for the white-box attack and the black-box attack under different channel conditions. For each channel condition, the eavesdropped image by the white-box attack is displayed on the left and the one obtained by the black-box attack is on the right.

of the images received by Bob, the eavesdropped images are visually recognizable and their privacy is compromised, which confirms the effectiveness of MIEA and reveals the risk of privacy leaks in current semantic communication.

TABLE I  
The average SSIM and PSNR of the eavesdropped images under different channel conditions

	Main Channel SNR		
	0dB	10dB	20dB
Reconstructed Images by Bob	30.02dB / 0.70	32.28dB / 0.79	33.28dB / 0.80
EC <sup>1</sup> SNR 0dB	17.44dB / 0.31	16.69dB / 0.30	18.58dB / 0.32
	21.65dB / 0.46	22.24dB / 0.46	22.85dB / 0.46
EC SNR 10dB	18.58dB / 0.40	17.80dB / 0.40	20.56dB / 0.43
	23.50dB / 0.53	23.76dB / 0.54	24.33dB / 0.56
EC SNR 20dB	18.74dB / 0.43	17.94dB / 0.42	20.84dB / 0.46
	23.88dB / 0.57	23.98dB / 0.59	24.41dB / 0.61

<sup>1</sup> EC refers to the eavesdropper channel.

### C. Evaluation of the Proposed Defense Method

Next, we evaluate the proposed defense method by repeating the evaluation of MIEA in section ?? with the defense method implemented on  $\mathbf{y}_f$ .

Fig. ?? visualizes the eavesdropped images by the two attack types after applying the proposed defense method, using the same individual as in Fig. ?. It can be observed that the eavesdropped images are visually unrecognizable, demonstrating the effectiveness of the proposed defense method in preventing Eve from eavesdropping on raw images. We can also see that the contour of the female in the white-box attack is less obvious than that in the black-box attack, suggesting that the defense against the white-box attack is superior to that against the black-box attack. This is because that Eve has no prior knowledge of the defense method when performing the white-box attack, whereas  $f^{-1}(\cdot)$  used in the black-box attack has learned some knowledge of the defense method from the training samples.

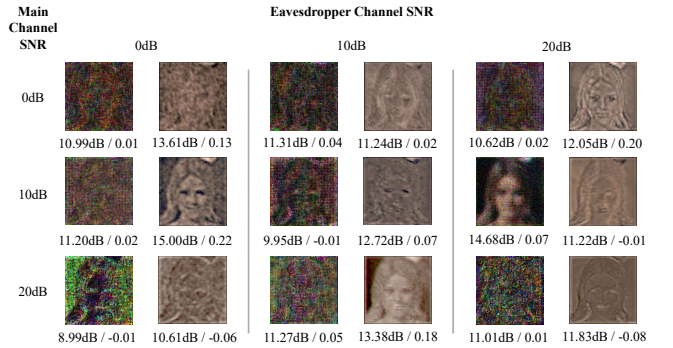


Fig. 4. Visualization of MIEA for both attacks under different channel conditions after applying the proposed method.

TABLE II  
The average SSIM and PSNR of the eavesdropped images by MIEA after defense

EC	MC <sup>1</sup>	0dB	10dB	20dB
0dB		8.03dB / 0.02	8.74dB / 0.02	6.94dB / 0.00
		11.36dB / 0.11	11.41dB / 0.07	12.51dB / 0.07
10dB		8.55dB / 0.04	7.70dB / -0.01	9.07dB / 0.05
		11.72dB / 0.16	11.34dB / 0.11	13.22dB / 0.13
20dB		8.02dB / 0.03	13.31dB / 0.09	8.39dB / 0.02
		12.59dB / 0.21	11.55dB / 0.10	11.54dB / 0.07

<sup>1</sup> MC refers to the main channel.

In addition, we provide the average SSIM and PSNR of the eavesdropped images for both attacks in Table. ?. For a given SNR of the main channel, the SSIM and PSNR do not increase as the SNR of the eavesdropper channel increases because different  $P$  and  $S$  are used for different transmitted features. The average SSIM and PSNR of the eavesdropped images by the black-box attack are larger than those by the white-box attack, which is consistent with the observation from Fig. ?. Overall, the average SSIM and PSNR are relatively small, which



indicates the effectiveness of the proposed defense method in preventing Eve from obtaining meaningful information from the eavesdropped signal.

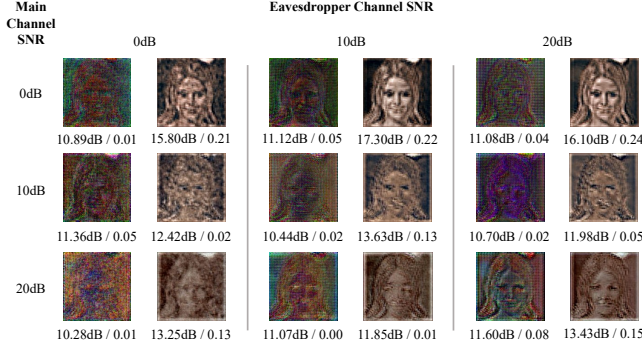


Fig. 5. Visualization of MIEA for both attacks under different channel conditions after applying only the random permutation.

TABLE III

The average SSIM and PSNR of the eavesdropped images when applying only random permutation

EC \ MC	0dB	10dB	20dB
0dB	8.05dB / 0.02 14.00dB / 0.22	8.73dB / 0.06 11.43dB / 0.10	8.50dB / 0.02 13.25dB / 0.13
10dB	8.32dB / 0.06 14.55dB / 0.26	8.02dB / 0.03 12.97dB / 0.19	8.87dB / 0.01 13.19dB / 0.12
20dB	8.36dB / 0.04 14.97dB / 0.28	8.05dB / 0.02 12.83dB / 0.18	8.77dB / 0.07 12.83dB / 0.13

Next, we conduct an ablation study to further validate our proposed method. Fig. ?? and Table. ?? demonstrate the eavesdropped images and the average SSIM and PSNR for both attacks by applying only the random permutation. As shown in Fig. ??, when only the random substitution is applied, the white-box attack can be effectively defended, while the black-box attack can still reconstruct visually recognizable images for some  $P$  and  $S$ . Moreover, the average SSIM and PSNR for the black-box attack in Table. ?? are larger than those in Table. ??, which demonstrates that only random permutation is insufficient for defending against MIEA.

Fig. ?? and Table. ?? show the related results for both attacks by applying only the random substitution. For both attacks, most of the eavesdropped images are visually recognizable, indicating that the attacker can still obtain sensitive information from the transmitted symbols, even though some of the semantic features have been substituted. The average SSIM and PSNR in Table. ?? are larger than those in Table. ??, which means that the random permutation is more effective than random substitution in defending against MIEA. From the ablation study, we can observe that the proposed defense method outperforms both the random-permutation-based and random-substitution-based defense methods, demonstrating that

both permutation and substitution are essential for the effectiveness of the proposed defense method.

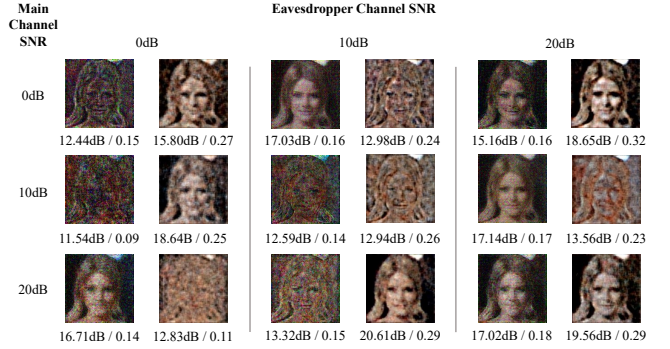


Fig. 6. Visualization of MIEA for both attacks under different channel conditions after applying only the random substitution.

TABLE IV

The average SSIM and PSNR of the eavesdropped images when applying only random substitution

EC \ MC	0dB	10dB	20dB
0dB	8.99dB / 0.14 16.06dB / 0.28	8.91dB / 0.10 15.00dB / 0.25	15.80dB / 0.16 14.70dB / 0.20
10dB	15.62dB / 0.19 14.68dB / 0.26	9.16dB / 0.14 14.49dB / 0.26	10.43dB / 0.16 15.80dB / 0.27
20dB	12.68dB / 0.18 15.78dB / 0.30	15.35dB / 0.18 14.53dB / 0.27	16.13dB / 0.21 17.29dB / 0.28

## V. Conclusion

In this paper, we propose MIEA to expose privacy risks in semantic communication. MIEA enables an attacker to eavesdrop on the transmitted symbols through an eavesdropper channel and reconstruct the raw message by inverting the DL model employed in the semantic communication system. We consider MIEA under the white-box attack and the black-box attack and propose a novel defense method based on random permutation and substitution to defend against both types of attack. In our evaluation, we first examine MIEA for both attacks under various channel conditions. We then conduct experiments and an ablation study to demonstrate the effectiveness of our proposed defense method.