

The Model Inversion Eavesdropping Attack in Semantic Communication Systems

Yuhao Chen, Qianqian Yang[†], Zhiguo Shi, Jining 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

{fsechenyh, qianqianyang20[†], shizg, qjm}@zju.edu.cn

Abstract—In recent years, semantic communication has been a popular research topic for its superiority in communication efficiency. This semantic communication relies on deep learning to extract information from raw messages, which is vulnerable to attacks targeting deep learning models. In this paper, we introduce the model inversion-based dropping attack (MIEA) to reveal the risk of privacy leaks in the semantic communication system. In MIEA, the attacker MIEA drops the signals 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 sample mutation 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-based defense method in preventing MIEA.

I.I. INTRODUCTION

Recently, semantic communication has been widely believed to be the next generation technology solution for 6G (6th generation) of wireless networks because of its high communication efficiency [2]. Compared with the current research, the current research which focuses on transmitting grouped bit sequences of the message [3] [4] [5] of the communication system, [6] [7] [8] [9] [10] [11] [12] [13] [14] [15] [16] [17] [18] [19] [20] [21] [22] [23] [24] [25] [26] [27] [28] [29] [30] [31] [32] [33] [34] [35] [36] [37] [38] [39] [40] [41] [42] [43] [44] [45] [46] [47] [48] [49] [50] [51] [52] [53] [54] [55] [56] [57] [58] [59] [60] [61] [62] [63] [64] [65] [66] [67] [68] [69] [70] [71] [72] [73] [74] [75] [76] [77] [78] [79] [80] [81] [82] [83] [84] [85] [86] [87] [88] [89] [90] [91] [92] [93] [94] [95] [96] [97] [98] [99] [100] [101] [102] [103] [104] [105] [106] [107] [108] [109] [110] [111] [112] [113] [114] [115] [116] [117] [118] [119] [120] [121] [122] [123] [124] [125] [126] [127] [128] [129] [130] [131] [132] [133] [134] [135] [136] [137] [138] [139] [140] [141] [142] [143] [144] [145] [146] [147] [148] [149] [150] [151] [152] [153] [154] [155] [156] [157] [158] [159] [160] [161] [162] [163] [164] [165] [166] [167] [168] [169] [170] [171] [172] [173] [174] [175] [176] [177] [178] [179] [180] [181] [182] [183] [184] [185] [186] [187] [188] [189] [190] [191] [192] [193] [194] [195] [196] [197] [198] [199] [200] [201] [202] [203] [204] [205] [206] [207] [208] [209] [210] [211] [212] [213] [214] [215] [216] [217] [218] [219] [220] [221] [222] [223] [224] [225] [226] [227] [228] [229] [230] [231] [232] [233] [234] [235] [236] [237] [238] [239] [240] [241] [242] [243] [244] [245] [246] [247] [248] [249] [250] [251] [252] [253] [254] [255] [256] [257] [258] [259] [260] [261] [262] [263] [264] [265] [266] [267] [268] [269] [270] [271] [272] [273] [274] [275] [276] [277] [278] [279] [280] [281] [282] [283] [284] [285] [286] [287] [288] [289] [290] [291] [292] [293] [294] [295] [296] [297] [298] [299] [300] [301] [302] [303] [304] [305] [306] [307] [308] [309] [310] [311] [312] [313] [314] [315] [316] [317] [318] [319] [320] [321] [322] [323] [324] [325] [326] [327] [328] [329] [330] [331] [332] [333] [334] [335] [336] [337] [338] [339] [340] [341] [342] [343] [344] [345] [346] [347] [348] [349] [350] [351] [352] [353] [354] [355] [356] [357] [358] [359] [360] [361] [362] [363] [364] [365] [366] [367] [368] [369] [370] [371] [372] [373] [374] [375] [376] [377] [378] [379] [380] [381] [382] [383] [384] [385] [386] [387] [388] [389] [390] [391] [392] [393] [394] [395] [396] [397] [398] [399] [400] [401] [402] [403] [404] [405] [406] [407] [408] [409] [410] [411] [412] [413] [414] [415] [416] [417] [418] [419] [420] [421] [422] [423] [424] [425] [426] [427] [428] [429] [430] [431] [432] [433] [434] [435] [436] [437] [438] [439] [440] [441] [442] [443] [444] [445] [446] [447] [448] [449] [450] [451] [452] [453] [454] [455] [456] [457] [458] [459] [460] [461] [462] [463] [464] [465] [466] [467] [468] [469] [470] [471] [472] [473] [474] [475] [476] [477] [478] [479] [480] [481] [482] [483] [484] [485] [486] [487] [488] [489] [490] [491] [492] [493] [494] [495] [496] [497] [498] [499] [500] [501] [502] [503] [504] [505] [506] [507] [508] [509] [510] [511] [512] [513] [514] [515] [516] [517] [518] [519] [520] [521] [522] [523] [524] [525] [526] [527] [528] [529] [530] [531] [532] [533] [534] [535] [536] [537] [538] [539] [540] [541] [542] [543] [544] [545] [546] [547] [548] [549] [550] [551] [552] [553] [554] [555] [556] [557] [558] [559] [560] [561] [562] [563] [564] [565] [566] [567] [568] [569] [570] [571] [572] [573] [574] [575] [576] [577] [578] [579] [580] [581] [582] [583] [584] [585] [586] [587] [588] [589] [590] [591] [592] [593] [594] [595] [596] [597] [598] [599] [600] [601] [602] [603] [604] [605] [606] [607] [608] [609] [610] [611] [612] [613] [614] [615] [616] [617] [618] [619] [620] [621] [622] [623] [624] [625] [626] [627] [628] [629] [630] [631] [632] [633] [634] [635] [636] [637] [638] [639] [640] [641] [642] [643] [644] [645] [646] [647] [648] [649] [650] [651] [652] [653] [654] [655] [656] [657] [658] [659] [660] [661] [662] [663] [664] [665] [666] [667] [668] [669] [670] [671] [672] [673] [674] [675] [676] [677] [678] [679] [680] [681] [682] [683] [684] [685] [686] [687] [688] [689] [690] [691] [692] [693] [694] [695] [696] [697] [698] [699] [700] [701] [702] [703] [704] [705] [706] [707] [708] [709] [710] [711] [712] [713] [714] [715] [716] [717] [718] [719] [720] [721] [722] [723] [724] [725] [726] [727] [728] [729] [730] [731] [732] [733] [734] [735] [736] [737] [738] [739] [740] [741] [742] [743] [744] [745] [746] [747] [748] [749] [750] [751] [752] [753] [754] [755] [756] [757] [758] [759] [760] [761] [762] [763] [764] [765] [766] [767] [768] [769] [770] [771] [772] [773] [774] [775] [776] [777] [778] [779] [780] [781] [782] [783] [784] [785] [786] [787] [788] [789] [790] [791] [792] [793] [794] [795] [796] [797] [798] [799] [800] [801] [802] [803] [804] [805] [806] [807] [808] [809] [810] [811] [812] [813] [814] [815] [816] [817] [818] [819] [820] [821] [822] [823] [824] [825] [826] [827] [828] [829] [83

consumable reliable transmission while simultaneously being transmittable to determine level of confidentiality protection. However, this is a sensitive communication system to provide compatibility and level of privacy protection. However, these reveal more privacy information systems. Secondly, deep learning based semantic communication is highly vulnerable to attacks targeting DLI models. For example, its have been conducted on attacks on the DLI model a review of available attacks regarding DLI models. Extensive studies being transmitted and dropped by a malicious attacker, the attacker can reconstruct the received message by utilizing the DL based data techniques. The attacker can also add perturbation in the transmitted data, attacking the semantic communication system to make it incoherent DL based attack techniques. For example, Sanyal *et al.* [9] proposed a multi-domain based attack causing the semantic communication system to make incorrect classifications, which is a novel idea. By introducing Sanyal input in a proposed semantic domain via a proposed semantic data poisoning attack, system causes the receiver to receive false, which is achieved by introducing noise at the receiver with the semantic feature. With a proposed image with a proposed attack, a black attack is performed by minimizing the difference between the semantic features of the targeted message and the received message. In this paper, we propose a semantic image with a proposed semantic communication system and by introducing the model diversion based dropping attack (MIEA) for semantic communication, where an attacker intercepts the transmitted symbols and attempts to reconstruct the original message from them by inverting the DLI systems used at the transmitter. We perform MIEA in both the white (MIEA_w) and black-box settings. The attacker has knowledge of the DLI models in the white-box settings but not in the black-box settings. To defend against MIEA, we propose a defense method based on a deep learning model and substitution. Evaluations demonstrate that the MIEA attack works under different channel conditions, in different values of the signal-to-noise ratio (SNR), which have a high risk of being attacked in MIEA in a transmission. A detailed result is based on the effectiveness of our proposed defense method. Evaluations demonstrate that the MIEA attack works in this paper is organized as follows. Section 2, introduces the basic of the signal-to-noise ratio (SNR) on which sections 2, we present the proposed MIEA under both the white-box and black-box settings 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. In section ??, we evaluate the effectiveness of the proposed MIEA and the proposed defense method. Section ?? concludes our work.

The transmitter of the semantic communication system consists of a semantic encoder and a channel encoder. The semantic encoder extracts the semantic features from the raw image, while the channel encoder maps the features into the transmitted symbols.

We consider a semantic communication system that 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 from the raw image, while the channel encoder maps the features into the transmitted symbols. We consider a semantic communication system that 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 from the raw image, while the channel encoder maps the features into the transmitted symbols.

Before transmission, y_f is reshaped into the transmitted symbols $y \in \mathbb{R}^{N \times 2}$, where $N = \frac{h \times w \times c}{2}$ and H and G are the real parts and imaginary parts of the signal to be transmitted, respectively. y is then transmitted over a wireless channel, such as multi-path propagation, fading and interference, while we denote as the main channel to distinguish from the channel used by the attacker. The received signal \hat{y} at the receiver side can be characterized by decoder and a semantic decoder. The receiver first reshapes \hat{y} back to the transmitted features \hat{y}_f . Then the channel decoder maps \hat{y}_f back to the semantic features \hat{z} . The semantic decoder then reconstructs the image \hat{x} from \hat{z} . We jointly train the semantic encoder, channel encoder, semantic decoder and channel decoder using the following loss function:

The receiver of the semantic communication system consists of a channel decoder and a semantic decoder. The receiver first reshapes \hat{y} back to the transmitted features \hat{y}_f . Then the channel decoder maps \hat{y}_f back to the semantic features \hat{z} . The semantic decoder then reconstructs the image \hat{x} from \hat{z} . We jointly train the semantic encoder, channel encoder, semantic decoder and channel decoder using the following loss function:

The first term in (??) computes the mean square error (MSE) between x and \hat{x} . The second term $T(x)$ is the total variation [?] that measures the smoothness of the reconstructed image \hat{x} , where $\hat{x}_{i,j}$ denotes the pixel value at the position (i, j) and β controls the smoothness of the image, with larger β being more piecewise-smooth. The hyper-parameter λ balances the two terms. In our work, we choose $\beta = 1$ and $\lambda = 1$.

Next, we introduce how an attacker eavesdrops the transmitted signal under the semantic communication system. We consider a scenario where an attacker Eve intercepts the transmitted signal 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

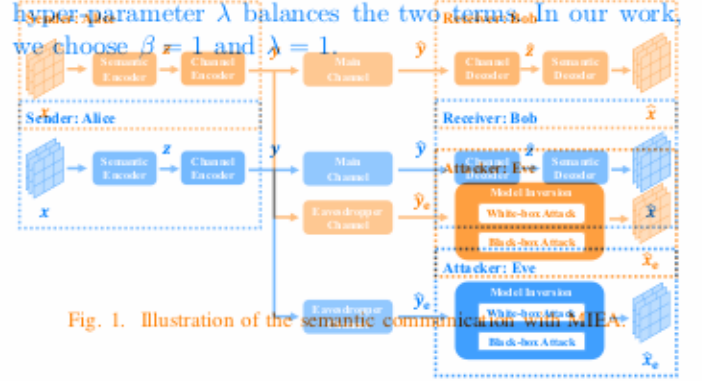


Fig. 1. Illustration of the semantic communication with MIEA.

it can easily be captured by any unauthorized receiver. Assume that there exists an attacker Eve who intercepts 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

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 y is transmitted over the wireless channel, it can easily be captured by eavesdropping y_e . Eve is able to reconstruct the image, any unauthorized receiver. Assume that there exists an attacker Eve who intercepts 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

Similarly, H_{ev} represents the eavesdropper channel matrix and n_e is a zero-mean additive white Gaussian noise. After eavesdropping y_e , Eve is able to reconstruct the image, any unauthorized receiver. Assume that there exists an attacker Eve who intercepts 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

III. THE PROPOSED MIEA AND ITS DEFENSE

In this section, we first elaborate the idea of MIEA. To reconstruct \hat{x} , 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. Note that to avoid confusion, we use the received image to denote \hat{x} received by Bob and the eavesdropped image to denote \hat{x}_e eavesdropped by Eve.

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 IPCEA.

In this section, we first elaborate the idea of MIEA. To reconstruct \hat{x} , 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. Note that to avoid confusion, we use the received image to denote \hat{x} received by Bob and the eavesdropped image to denote \hat{x}_e eavesdropped by Eve.

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 IPCEA. To reconstruct \hat{x} , 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. Note that to avoid confusion, we use the received image to denote \hat{x} received by Bob and the eavesdropped image to denote \hat{x}_e eavesdropped by 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 transmission 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 [1] to transmit images from the CelebA dataset [8], while the channel encoder and decoder have one convolutional layer. The semantic encoder and decoder each have four convolutional layers. 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 three distinct channel conditions. 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.

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 [1], the SNR value determines the standard deviation of γ when the transmission power is normalized to 1. We train the CelebA dataset with a batch size of 28 using Adam (0dB, 10dB and 20dB) a learning rate of 10^{-3} distinct DeepJSCC models. As stated in [2], the SNR value determines the standard deviation of γ when the signal is normalized to 1. We train the CelebA dataset with a batch size of 28 using Adam (0dB, 10dB and 20dB) a learning rate of 10^{-3} distinct DeepJSCC models.

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) [21], on the CelebA dataset. SSIM and PSNR indicate better quality. We employ Adam [2] as the optimizer with a learning rate of 10^{-3} .

B. Evaluation of MIEA To measure the image quality, we use two metrics, i.e., the SSIM and PSNR. 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 [2] 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 \mathcal{X} and obtain its corresponding transmitted symbols \mathcal{Y} . Similarly, we use Adam as the optimizer and set the learning rate to 10^{-3} .

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 [2] 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 \mathcal{X} and obtain its corresponding transmitted symbols \mathcal{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 image transmitted by Alice and the second column shows the images received by Bob. For the black-box attack, we use the images received by Bob consisting of different upsampling layers and convolution layers. Then two training sets, (i) increasing the SNR optimization problem quality, which is related to high CelebA images SSIM and PSNR, and (ii) decreasing the SNR optimization problem quality, which is related to low CelebA images SSIM and PSNR, are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used.

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 image transmitted by Alice and the second column shows the images received by Bob. For the black-box attack, we use the images received by Bob consisting of different upsampling layers and convolution layers. Then two training sets, (i) increasing the SNR optimization problem quality, which is related to high CelebA images SSIM and PSNR, and (ii) decreasing the SNR optimization problem quality, which is related to low CelebA images SSIM and PSNR, are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used. Additionally, higher SNRs symbols \mathcal{Y} of CelebA dataset are used.

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 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.

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.

Reconstructed Images by Bob	Main Channel SNR		
	0dB	10dB	20dB
EC SNR 0dB	30.92dB / 0.70	32.28dB / 0.79	33.28dB / 0.80
EC SNR 10dB	17.44dB / 0.31	18.69dB / 0.30	18.58dB / 0.32
EC SNR 20dB	18.58dB / 0.40	17.80dB / 0.40	20.56dB / 0.43
EC SNR 0dB	23.50dB / 0.53	23.76dB / 0.54	24.33dB / 0.56
EC SNR 10dB	18.74dB / 0.43	17.94dB / 0.42	20.84dB / 0.46
EC SNR 20dB	23.88dB / 0.57	23.98dB / 0.59	24.41dB / 0.61

TABLE I
The average SSIM and PSNR of the eavesdropped images under different channel conditions. EC refers to the eavesdropper channel.

Reconstructed Images by Bob	Main Channel SNR		
	0dB	10dB	20dB
EC SNR 0dB	30.92dB / 0.70	32.28dB / 0.79	33.28dB / 0.80
EC SNR 10dB	17.44dB / 0.31	18.69dB / 0.30	18.58dB / 0.32
EC SNR 20dB	18.58dB / 0.40	17.80dB / 0.40	20.56dB / 0.43
EC SNR 0dB	23.50dB / 0.53	23.76dB / 0.54	24.33dB / 0.56
EC SNR 10dB	18.74dB / 0.43	17.94dB / 0.42	20.84dB / 0.46
EC SNR 20dB	23.88dB / 0.57	23.98dB / 0.59	24.41dB / 0.61

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

Reconstructed Images by Bob	Main Channel SNR		
	0dB	10dB	20dB
EC SNR 0dB	30.92dB / 0.70	32.28dB / 0.79	33.28dB / 0.80
EC SNR 10dB	17.44dB / 0.31	18.69dB / 0.30	18.58dB / 0.32
EC SNR 20dB	18.58dB / 0.40	17.80dB / 0.40	20.56dB / 0.43
EC SNR 0dB	23.50dB / 0.53	23.76dB / 0.54	24.33dB / 0.56
EC SNR 10dB	18.74dB / 0.43	17.94dB / 0.42	20.84dB / 0.46
EC SNR 20dB	23.88dB / 0.57	23.98dB / 0.59	24.41dB / 0.61

Next, we evaluate the proposed defense method by repeating the evaluation of MIEA in section 3.4 with the defense method implemented on the eavesdropped images by the white-box attack and the black-box attack. As in Fig. 3, 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 distinct 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, while, as in Fig. 3, the black-box attack has learned some knowledge of the defense method from the training samples. The effectiveness of the proposed defense method in preventing Eve from eavesdropping on raw images is further confirmed by the average SSIM and PSNR of the eavesdropped images for both attacks in Table 2. For a given SNR of the main channel, the SSIM and PSNR of the eavesdropped images by the black-box attack are higher than those by the white-box attack, which is consistent with the observation from Fig. 3. Overall, the average SSIM and PSNR are relatively small, which indicates the effectiveness of

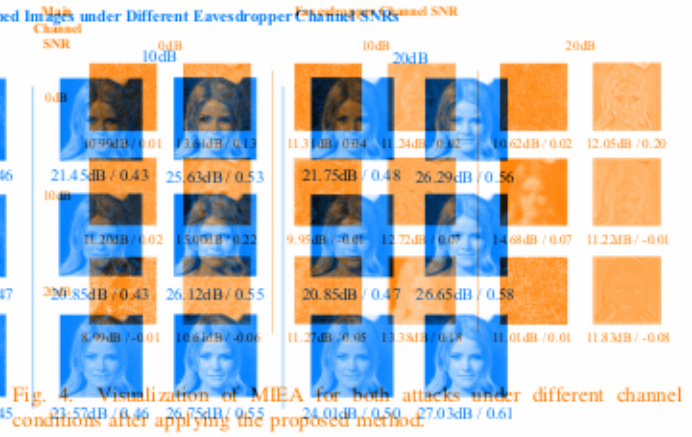


Fig. 4. Visualization of MIEA for both attacks under different channel conditions after applying the proposed method.

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

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

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

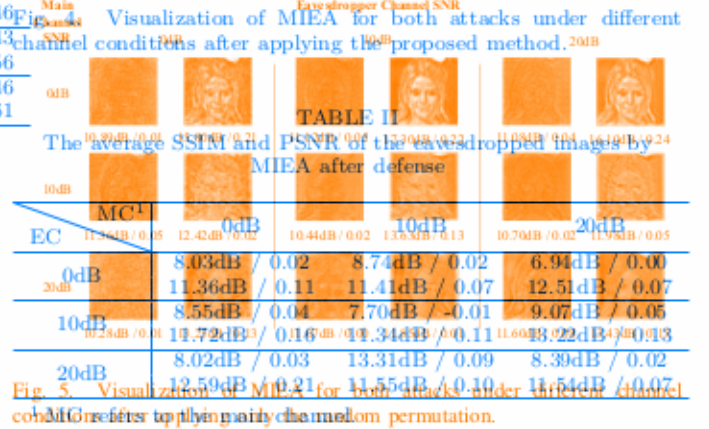


Fig. 5. Visualization of MIEA for both attacks under different channel conditions after applying the proposed method.

Next, we conduct an ablation study to further validate our proposed method. Fig. 7 and Table 3 show the SSIM and PSNR of the eavesdropped images for both attacks and PSNR for both attacks given SNR of the main channel. The SSIM and PSNR of the eavesdropped images are very low, which indicates that the white-box attack can be effectively defended, while the black-box attack is still not so secure. The average SSIM and PSNR of the eavesdropped images by the black-box attack are higher than those by the white-box attack, which is consistent with the observation from Fig. 3. Overall, the average SSIM and PSNR are relatively small, which indicates the effectiveness of

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

EC \ MC		0dB	10dB	20dB
Main Channel SNR	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
	10dB	11.54dB / 0.09 16.71dB / 0.14	12.44dB / 0.15 12.83dB / 0.11	12.99dB / 0.14 13.32dB / 0.15

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.

EC \ MC		0dB	10dB	20dB
Main Channel SNR	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
	10dB	12.44dB / 0.15 16.71dB / 0.14	12.99dB / 0.14 13.32dB / 0.15	13.56dB / 0.23 19.56dB / 0.29

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 .

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

EC \ MC		0dB	10dB	20dB
Main Channel SNR	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
	10dB	17.03dB / 0.16 17.02dB / 0.18	12.99dB / 0.14 13.32dB / 0.15	13.56dB / 0.23 19.56dB / 0.29

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.

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.

both depend 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.

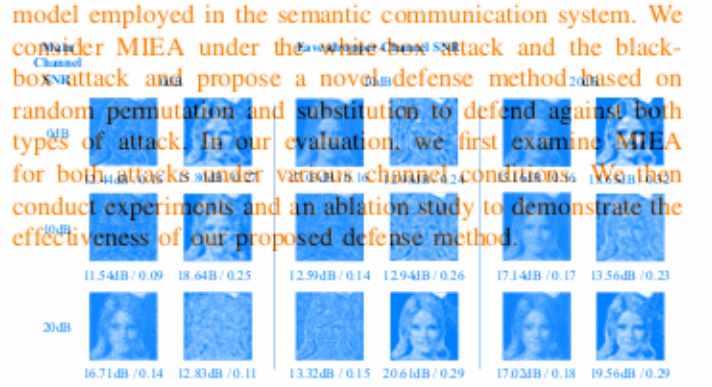


Fig. 6. Visualization of MIEA for both attacks under different channel conditions after applying only the random substitution.

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

EC \ MC		0dB	10dB	20dB
Main Channel SNR	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
	10dB	17.03dB / 0.16 17.02dB / 0.18	12.99dB / 0.14 13.32dB / 0.15	13.56dB / 0.23 19.56dB / 0.29

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.