

# Secure Deep-JSCC Against Multiple Eavesdroppers

Seyyed Amirhossein Ameli Kalkhoran<sup>\*†</sup>, Mehdi Letafati<sup>\*†</sup>, Ecenaz Erdemir<sup>§</sup>, Babak Hossein Khalaj<sup>†</sup>,  
Hamid Behroozi<sup>†</sup>, and Deniz Gündüz<sup>§</sup>

<sup>†</sup> Electrical Engineering Department, Sharif University of Technology, Tehran, Iran

<sup>†</sup> Computer Engineering Department, Sharif University of Technology, Tehran, Iran

<sup>§</sup> Department of Electrical and Electronic Engineering, Imperial College London, UK

Emails: †{mletafati@ee., behroozi@, khalaj@}sharif.edu; †ameli@ce.sharif.edu; §{e.erdemir17, d.gunduz}@imperial.ac.uk

**Abstract**—In this paper, a generalization of deep learning-aided joint source channel coding (Deep-JSCC) approach to secure communications is studied. We propose an end-to-end (E2E) learning-based approach for secure communication against multiple eavesdroppers over complex-valued fading channels. Both scenarios of colluding and non-colluding eavesdroppers are studied. For the colluding strategy, eavesdroppers share their logits to collaboratively infer private attributes based on ensemble learning method, while for the non-colluding setup they act alone. The goal is to prevent eavesdroppers from inferring private (sensitive) information about the transmitted images, while delivering the images to a legitimate receiver with minimum distortion. By generalizing the ideas of privacy funnel and wiretap channel coding, the trade-off between the image recovery at the legitimate node and the information leakage to the eavesdroppers is characterized. To solve this secrecy funnel framework, we implement deep neural networks (DNNs) to realize a data-driven secure communication scheme, without relying on a specific data distribution. Simulations over CIFAR-10 dataset verifies the secrecy-utility trade-off. Adversarial accuracy of eavesdroppers are also studied over Rayleigh fading, Nakagami- $m$ , and AWGN channels to verify the generalization of the proposed scheme. Our experiments show that employing the proposed secure neural encoding can decrease the adversarial accuracy by 28%.

**Index Terms**—Secure Deep-JSCC, data-driven security, secrecy-utility trade-off, secure image transmission.

## I. Introduction

Driven by the growing interest in semantic communication systems [?], intelligent transmission of multimedia content has received much attention because of its various applications in augmented/virtual reality (AR/VR), Metaverse [?], and surveillance systems [?], [?]. The adoption and success of such services rely highly on the security of the delivered contents—communication systems should understand the desired “level of security” and intelligently adapt the transmission scheme accordingly [?], [?].

Connected intelligence is foreseen as the most significant driving force in the sixth generation (6G) of wireless communications. To this end, artificial intelligence and machine learning (AI/ML) algorithms are envisioned to be widely incorporated into 6G networks, realizing an “AI-native” air interface. Nevertheless, security issues at the wireless edge of 6G networks are still identified as open challenges [?]. The air interface of 6G systems encounters

ever-rising attacks, such as eavesdropping, spoofing [?], and man-in-the-middle [?].

Recently, a considerable number of research has been dedicated to the utilization of deep learning (DL) techniques to optimize the performance of wireless systems, thanks to their outstanding performance and generalization capabilities [?], [?], [?]. In the context of wireless security, autoencoders (composed of linear layers) are exploited in [?] over the additive white Gaussian noise (AWGN) wiretap channel. To tackle the trade-off between the data rate and security, a weighted sum of block error rate and information leakage is used as the loss function (LF) for neural wiretap code design. The data fed into the autoencoder is combined with additional non-informative random bits to confuse the eavesdropper; while, this also reduces the communication rate. Notably, most of the previous works, i.e., [?], [?], [?], focus on learning-aided secure channel coding, rather than taking into account the end-to-end (E2E) performance of secure communications. The content of the transmitted data is not addressed in these works and the entire bit-stream is equally treated as the secret information to be protected against an eavesdropper.

The E2E communication of images from a source node to a legitimate destination can be considered joint source channel coding (JSCC) problem. DL-JSCC design, Deep-JSCC, has received significant attention thanks to its superior performance, particularly its better resilience to accurate channel state information [?]. However, JSCC different from separate source and channel coding, the channel codeword is correlated with the underlying source signal. This can create vulnerabilities in terms of leakage to eavesdroppers, despite providing robustness against channel noise. Inspired by [?] and [?], we provide a generalization of the Deep-JSCC approach to secure communication problems against multiple eavesdroppers. In this regard, [?] proposes a generative adversarial network (GAN)-inspired secure neural encoder-decoder pair over an AWGN wiretap channel against one eavesdropper. The authors in [?] propose a variational autoencoder (VAE)-based approach for Deep-JSCC design over binary symmetric channels, again considering a single eavesdropper.

In this paper, we consider E2E learning-based secure

<sup>\*</sup>Equal contribution







ability distribution of the original and the reconstructed image, taking into account the randomness in the input image and the channel. Since the true distribution  $p(\mathbf{u})$  is often unknown, we estimate the expected distortion measure using samples  $\mathbf{u}_j$  from an available dataset  $\mathcal{D}_u$  by computing  $\mathbb{E}_{p(\mathbf{u}, \hat{\mathbf{u}})}[d(\mathbf{u}, f_{\Omega_B}(\mathbf{y}))] \approx \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} d(\mathbf{u}, \hat{\mathbf{u}})$ , where  $N_u \triangleq |\mathcal{D}_u|$ . It is assumed that we know the sensitive attribute in which the eavesdroppers are interested, as well as their channel models. Both of these assumptions are common in the privacy [?], [?] and wiretap channel [?], [?], [?] literature. We do not need to know the instantaneous channel gains, but use their distributions to sample channel realizations during training.

**Training Procedure:** In order to train our system based on (??), we follow an iterative procedure. Intuitively, the network nodes are faced with a minimax game, i.e., the competition between legitimate autoencoder and the adversarial DNNs. Hence, the following strategy is run through our proposed E2E system: The encoder and decoder function of Alice and Bob should jointly minimize their LF, denoted by  $\mathcal{L}_{AB}$ :

$$\mathcal{L}_{AB} = \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} \left( d(\mathbf{u}, \hat{\mathbf{u}}) - \frac{1}{M} \sum_{m \in [M]} w_m H(q_{\Theta_{E,m}}(s(\mathbf{u}) | \mathbf{z}_m(\mathbf{u})), \varepsilon_s(\mathbf{u})) \right) \quad (5)$$

The training process of legitimate nodes can be further enhanced via employing adversarial likelihood compensation (ALC), which has been shown in [?] to be more effective in confusing an adversary than the one-hot encoding approach. The main idea is to make the posterior distribution of adversaries imitate a uniform distribution  $\bar{p}_L = [\frac{1}{L}, \dots, \frac{1}{L}]^T$ . Hence, Alice and Bob jointly maximize the uncertainty of adversarial predictions, resulting in the following loss function

$$\mathcal{L}_{AB}^{\text{ALC}} = \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} \left( d(\mathbf{u}, \hat{\mathbf{u}}) + \frac{1}{M} \sum_{m \in [M]} w_m H(q_{\Theta_{E,m}}(s(\mathbf{u}) | \mathbf{z}_m(\mathbf{u})), \bar{p}_L) \right) \quad (6)$$

The distortion measure we consider for our legitimate loss function  $\mathcal{L}_{AB}$  is a mixture of the average mean squared error (MSE), denoted by  $\Delta^{\text{MSE}}$ , and the structural similarity index (SSIM),  $\Delta^{\text{SSIM}}$ , between the input image  $\mathbf{u}$  and the recovered version  $\hat{\mathbf{u}}$  at the output of Bob's DNN. Therefore, we assume  $d(\cdot, \cdot)$  to be measured as follows

$$d(\mathbf{u}, \hat{\mathbf{u}}) = \Delta^{\text{MSE}}(\mathbf{u}, \hat{\mathbf{u}}) + \alpha \Delta^{\text{SSIM}}(\mathbf{u}, \hat{\mathbf{u}}), \quad (7)$$

where  $\Delta^{\text{MSE}}(\mathbf{u}, \hat{\mathbf{u}}) \triangleq \frac{1}{n} \|\mathbf{u} - \hat{\mathbf{u}}\|^2$ ,  $\Delta^{\text{SSIM}}(\mathbf{u}, \hat{\mathbf{u}}) \triangleq 1 - \text{SSIM}(\mathbf{u}, \hat{\mathbf{u}})$ , and  $\alpha$  is a tuning parameter representing the contribution of the SSIM metric. The SSIM measure

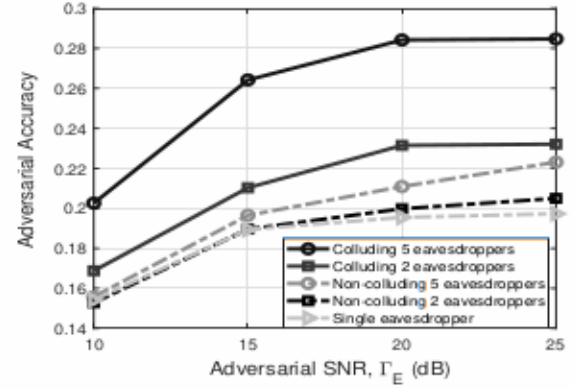


Fig. 4: Total adversarial accuracy over Rayleigh channels.

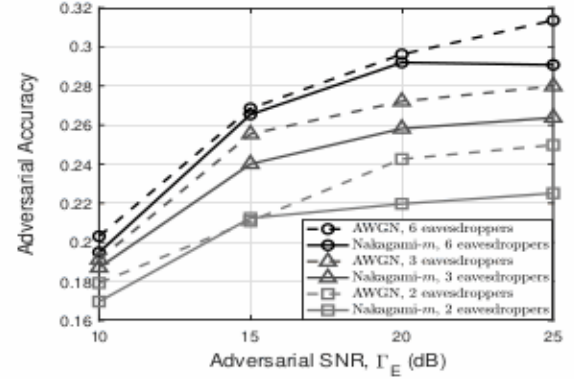


Fig. 5: Colluding adversarial accuracy over AWGN and Nakagami.

between two images  $I$  and  $K$  is defined as

$$\text{SSIM}(I, K) \triangleq \frac{(2\mu_I\mu_K + c_1)(2\sigma_{IK} + c_2)}{(\mu_I^2 + \mu_K^2 + c_1)(\sigma_I^2 + \sigma_K^2 + c_2)}, \quad (8)$$

where  $\mu_I$ ,  $\mu_K$ ,  $\sigma_I$ ,  $\sigma_K$ , and  $\sigma_{IK}$  are the local means, standard deviations, and cross-covariance for images  $I$  and  $K$ , while  $c_1$  and  $c_2$  are two adjustable constants [?]. The rationale behind the proposed distortion metric is that we not only aim to recover every pixel of images with minimum error (captured via the MSE measure), but also want to obtain a good-quality reconstruction from the human perception point-of-view.

Each step of the joint training of  $\mathcal{A}$ -to- $\mathcal{B}$  autoencoder is followed by a training step for the adversarial DNNs. Eavesdroppers aim to minimize the CE between their estimated likelihood  $q_{\Theta_{E,m}}(s | \mathbf{z}_m)$  and the ground-truth vector  $\varepsilon_s$  corresponding to  $S$ . Hence, the following LF is employed for training the DNN of Eve<sub>m</sub> for  $m \in [M]$ :

$$\mathcal{L}_{E,m} = \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} H(q_{\Theta_{E,m}}(s(\mathbf{u}) | \mathbf{z}_m(\mathbf{u})), \varepsilon_s(\mathbf{u})). \quad (9)$$

Note that the adversarial DNNs can be trained in parallel. Then for the case of colluding eavesdroppers, an additional step of “knowledge sharing” is performed. In this case, the adversaries share their individually-extracted logits, and a weighted sum of these logits is exploited for the inference of private attributes, where the logit weights are trained in the colluding framework.



#### IV. Evaluations

In this section, we evaluate the performance of the proposed scheme over both AWGN and complex fading (Rayleigh and Nakagami- $m$ ) communication channels. We address the generalization capability of our proposed scheme for different communication scenarios and over a wide range of signal-to-noise ratio (SNR) values, to highlight the data efficiency of our proposed learning-based security solution. We also address the secrecy-utility trade-off for the proposed learning-based approach. We evaluate our proposed secure framework using images with dimension  $32 \times 32 \times 3$  (height, width, channels) from CIFAR-10 dataset [?]. The dataset consists of 60000 colored images of size  $32 \times 32$  pixels. The training and evaluation sets are two completely separated sets of images, containing 50000 and 10000 images, respectively, associated with 10 classes. Adversaries wish to infer a common secret  $S$ , either individually (the non-colluding case), or by learning from the combination of the adversarial logits (the colluding setup). The common secret  $S$  here is considered as the class of CIFAR-10 images with  $|S| = L = 10$ . For simplicity, we consider a single weight in the LF, that is,  $w_m = w = 5, \forall m$ . We also set  $\alpha = 0.1$  and  $n_T = 4$  antennas. These parameters are set after conducting extensive experiments and training the DNNs with a wide range of values for  $w_m$  and  $\alpha$ , where we have omitted the results of fine-tuning step due to space limitations. Transmit SNRs of communication links are defined as  $\Gamma_B = 10 \log_{10} \frac{P}{\sigma_L^2}$  dB and  $\Gamma_E = 10 \log_{10} \frac{P}{\sigma_E^2}$  dB, representing the ratio of the average power of the channel input to the average noise power of legitimate  $\sigma_L^2$  and adversarial nodes  $\sigma_E^2$ , respectively. During training, we set  $\Gamma_B = 20$  dB and  $\Gamma_E = 15$  dB, respectively, while we test the performance over different values of channel SNRs during the inference. In addition, the bandwidth compression ratio is set to  $\frac{k}{n} = \frac{1}{3}$ . For the training, we sample channel realizations from the general case of complex Rayleigh fading model with average  $\Gamma_B$  and  $\Gamma_E$  values stated above. Nevertheless, during the inference phase we study the performance in different scenarios of AWGN and Nakagami- $m$  channels (for  $m = 3$ ). While we do assume known channel models in our simulations, which we use to generate samples from conditional channel distribution, we could easily drop this assumption if we had data collected from a particular channel with unknown statistics. DNN architectures are implemented using Python3 with Tensorflow.<sup>3</sup> The codes were run on Intel(R) Xeon(R) Silver 4114 CPU running at 2.20 GHz with GeForce RTX 2080 Ti GPU. To minimize the LFs, the widely-adopted Adam optimizer is chosen [?] with a learning rate of  $10^{-4}$ . We fix the number of training episodes to  $N_{\text{episode}} = 200$ , and the batch size to  $m = 128$ .

Figs. ?? and ?? show the adversarial accuracy of the proposed scheme vs. the SNR of adversarial links

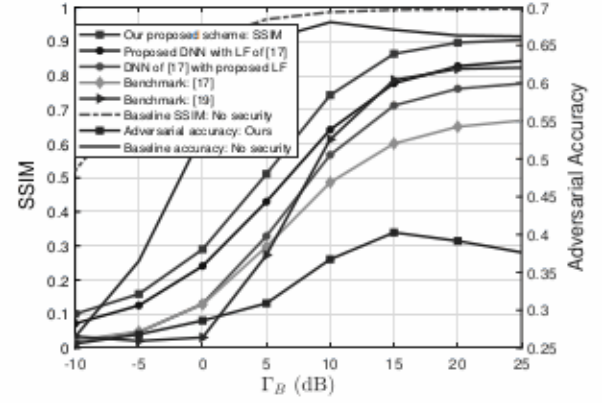


Fig. 6: Ablation study for  $M = 3$  non-colluding eavesdroppers.

( $\Gamma_E = \Gamma_B - 5$  dB) for both scenarios of colluding and non-colluding eavesdroppers. For the colluding setup, the total adversarial accuracy refers to the overall accuracy of adversaries in correctly finding the ground-truth  $\varepsilon_s$  from their aggregated logits, while for the non-colluding benchmarks, the mean accuracy across eavesdroppers is plotted for the sake of comparison. One can observe that increasing the number of eavesdroppers leads to higher accuracy for the adversaries, which is aligned with one's intuition. The increase in adversarial accuracy is more significant in the colluding case due to the collaboration and “knowledge sharing” among eavesdroppers through the ensemble learning process [?]. This actually helps them learn the secret more accurately. The figures also indicate that by increasing the quality of adversarial links, i.e., increasing  $\Gamma_E$ , the accuracy of adversaries increases by at most 10%. This is because higher SNR values result in having less-distorted (less noisy) observations at the eavesdroppers, resulting in more accurate estimations about the posterior adversarial distribution  $q_{\Theta_{E,m}}(s|z_m)$ . The amount of increase in the adversarial accuracy reduces with the increase in  $\Gamma_E$ , which highlights the limitation of eavesdroppers in the proposed secure scheme. Fig. ?? highlights the generalization capability of a fully learning-based framework extended to the AWGN and Nakagami- $m$  (3 channel) scenario. It shows that we can achieve almost similar results in these scenarios, other than Rayleigh fading, despite training the networks with a Rayleigh channel model.

Fig. ?? illustrates the data reconstruction performance at Bob and the total adversarial accuracy, when having  $M = 3$  non-colluding eavesdroppers. One can infer from the figure that our proposed system outperforms the benchmarks in terms of the reconstruction performance. Accordingly, 20% and 10% performance gain is achieved by our proposed scheme compared with [?] and [?], respectively. The ablation examinations conducted in this figure show that both the implemented DNNs and the proposed LFs for optimizing the framework contribute to the system's performance compared with other benchmarks. The figure also implies that increasing  $\Gamma_B$  results

<sup>3</sup><https://www.tensorflow.org/>

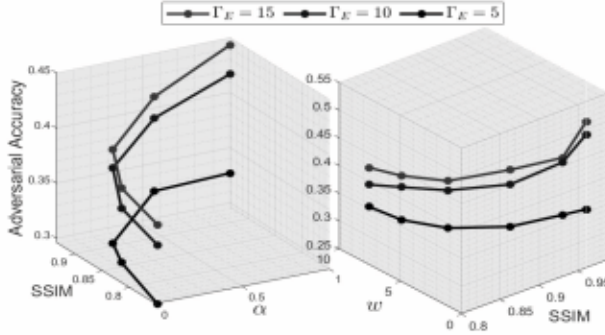


Fig. 7: Secrecy-utility trade-off.

in having higher SSIM values. This is because increasing  $\Gamma_B$  can result in less distorted observations at Bob, which facilitates the image reconstruction performance. Data efficiency and generalizability of our proposed scheme are also validated, since we have trained our DNNs with a fixed SNR  $\Gamma_B = 20$  dB, while the performance gain of our approach during inference holds for various SNRs. Furthermore, Fig. ?? highlights that if we ignore the eavesdroppers during the training of  $\mathcal{A}$ - $\mathcal{B}$  pair, and set  $w_m = 0, \forall m \in [M]$ , our proposed scheme can achieve almost perfect (SSIM = 1) data recovery. The impact of the eavesdroppers on Bob's performance can be studied in this figure as well, where having 3 eavesdroppers can impose 10% decrease in the reconstruction performance of Bob. Finally, to indicate the importance of our proposed adversary-aware scheme in terms of preventing leakage, we can observe from the figure that if we do not employ secure neural encoding (i.e., ignoring the eavesdroppers during the training of  $\mathcal{A}$ - $\mathcal{B}$  pair), the adversarial accuracy is increased by about 28%. This clearly highlights the importance of employing our proposed learning-based secure encoding scheme.

Fig. ?? studies the impact of tuning parameters  $\alpha$  and  $w$  on the adversarial performance of our proposed system. These hyper-parameters are the coefficients associated with utility and secrecy adjustment terms within our training LFs in (??) and (??), respectively. For this experiment, the adversarial performance is captured by investigating the accuracy of adversaries in correctly finding the ground-truth label  $\varepsilon_s$  (representing the sensitive information  $S$ ) among the labels of CIFAR-10 dataset. The figure indicates that by increasing  $\alpha$ , higher values of SSIM can be achieved, since more emphasis on the SSIM error  $\Delta^{\text{SSIM}}$  is put based on (??). However, the accuracy of adversaries in extracting the sensitive information also increases, which verifies the secrecy-utility trade-off. Notably, by increasing  $\alpha$  to values more than 0.1, a jump in the adversarial accuracy can be observed, which leads us choosing  $\alpha = 0.1$  for our network. Similarly, by increasing  $w$ , the emphasis goes toward the secrecy criteria introduced in (??), (??), and (??), which leads to the reduction in adversarial accuracy and achieved SSIM, verifying the secrecy-utility trade-off as well. Fig. ?? also shows that increasing  $\Gamma_E$  can

improve the adversarial accuracy in finding the sensitive data  $\varepsilon_s$ . Notably, the amount of increase in the adversarial accuracy reduces with the increase in  $\Gamma_E$  which highlights the limitation of adversarial nodes based on our proposed secure scheme.

## V. Conclusions and Future Directions

We proposed a learning-based approach for E2E secure image delivery against multiple eavesdroppers over AWGN and complex-valued fading channels. We considered both scenarios of colluding and non-colluding eavesdroppers over CIFAR-10 dataset. For the colluding strategy, eavesdroppers collaborate to infer private data from their observations (channel outputs), using ensemble learning, while for the non-colluding setup they act alone. Meanwhile, the legitimate parties aim to have a secure communication with minimum average distortion. Employing autoencoders, we proposed a secrecy funnel framework to achieve both secrecy and utility, where we also take into account the perceptual quality of image transmission within our LF. Evaluations validate the performance of our proposed scheme compared with existing benchmarks, while addressing the secrecy-utility trade-off. Our proposed system is also shown to be generalizable to a wide range of SNRs and different communication scenarios.

Future Directions: A problem to be studied is training the systems over real-time wireless channels. The challenge here is the relatively long coherence time compared to the rate at which data samples can be processed for training. Accordingly, only a few channel realizations are observed over every minibatch, which will be an important issue for training communication systems that are supposed to generalize well to a wide range of channel realizations.

## References

- [1] D. Gündüz, Z. Qin, I. Estella Aguerri, H. S. Dhillon, Z. Yang, A. Yener, K. K. Wong, and C.-B. Chae, "Beyond transmitting bits: Context, semantics, and task-oriented communications," arXiv:2207.09353, 2022.
- [2] M. Letafati and S. Sotiriou, "On the privacy and security of health data in the metaverse: An overview," *Adg.*, 2023, <https://doi.org/10.1016/j.adhoc.2023.103262>.
- [3] Q. Yang and D. Gündüz, "Coded caching and content delivery with heterogeneous distortion requirements," *IEEE Trans. Inf. Theory*, vol. 64, no. 6, pp. 4347–4364, Jun. 2018.
- [4] M. Jankowski, D. Gündüz, and K. Mikolajczyk, "Wireless image retrieval at the edge," *IEEE J. Sel. Areas Comm.*, vol. 39, no. 1, pp. 89–100, Jan. 2021.
- [5] M. Letafati, H. Behroozi, B. H. Khalaj, and E. A. Jorswieck, "Content-based medical image transmission against randomly-distributed passive eavesdroppers," *IEEE Int. Conf. Comm. Workshop (ICCW)*, Montreal, Canada, Jun. 2021, pp. 1–7.
- [6] —, "On learning-assisted content-based secure image transmission for delay-aware systems with randomly-distributed eavesdroppers," *IEEE Trans. Comm.*, vol. 70, no. 2, pp. 1139–1150, Feb. 2022.
- [7] A. Chorti, A. N. Barreto, S. Kopsell, M. Zoli, M. Chaffi, P. Sehier, G. Fettweis, and H. V. Poor, "Context-aware security for 6G wireless: The role of physical layer security," 2021, arXiv:2101.01536.