

# Secure Deep-JSCC Against Multiple Eavesdroppers

Seyyed Amirhossein Amirikhalil<sup>††</sup>, Mehdi Lotafati<sup>†\*</sup>, Erenaz 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: [f.molatafati@ee.ubt.ac.ir](mailto:f.molatafati@ee.ubt.ac.ir), [behroozi@sharif.edu](mailto:behroozi@sharif.edu), [khalaj@sharif.edu](mailto:khalaj@sharif.edu), [fameli@ee.sharif.edu](mailto:fameli@ee.sharif.edu), [c.erdemir17@gmail.com](mailto:c.erdemir17@gmail.com), [d.gunduz@imperial.ac.uk](mailto: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 a generalization of (E2E) learning-based approach for secure communication against multiple eavesdroppers over complex 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 their logits to collaboratively infer private attributes based on ensemble learning method, while for the non-colluding setup they infer (sensitive) information about the transmitted images with minimum distortion. By generalizing the idea of privacy funnel and wiretap channel coding the trade-off between the image reconstruction fidelity and the information leakage to the legitimate node and the information leakage to the eavesdroppers is characterized. To solve this *privacy 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 accuracy-utility trade-off. Adversarial accuracy of adversarial privacy in 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 [2], intelligent transmission of multimedia content has received much attention because of its various applications in augmented/virtual reality (AR/VR), Metaverse [3], and surveillance systems [4], [5]. 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 [4], [5].

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 [?].

ever-rising attacks, such as eavesdropping, spoofing [2], and recently, the considerable number of research has been dedicated to the utilization of deep learnings (DL) techniques to optimize the performance of wireless systems. DL thanks to their outstanding performance and generalizable capabilities [1], [2], [3]. In the context of wireless security, autoencoders (composed of encoder [4], layers [5]) are exploited in [2] to overcome the additive white Gaussian noise (AWGN) of wireless channels. To tackle the trade-off between the data rate and security in a (AWGN) sum of block error rate and information leakage is used as the loss function (LF) of a original winmap code design. The data fed into the autoencoder is combined with additional (LF)-informative winmap code design (the so-called dropper while, this also reduces the communication rate. Notably, none of the previous works i.e. the [2], [3], [4] focuses on training autoencoder in the channel condition rather than taking into account the channel work (E2E) performance [6] of secure communication. The content of the transmitted data is not addressed in these works and the (E2E) probability is equally treated as the secret information to be protected against data eavesdropper in the E2E communication of bit streams. A equal treatment as a legitimate destination becomes considered as joint source channel coding (JSCC) problem. DL-aided JSCC design, a.k.a. DeepJSCC has received significant attention thanks to its superior performance, particularly, its lack of dependence on joint source channel state information (JSCC) [7]. However, DL-JSCC design, separate DeepJSCC channel coding, this channel code word is correlated with the underlying source signal. This its lack of vulnerability in terms of leakage to eavesdropper, [8]. However, providing JSCC is robust against the eavesdropper is inspired by [9] and [10] yet provide a generalization of the DeepJSCC approach to secure communication. This problem against multiple eavesdroppers. In this regard, [2] droppers, a generative adversarial network (GAN) inspired secure inspirations by [11] decoder pair provide AWGN winmap code in the DeepJSCC eavesdropper. The authors in [11] propose a prohibition against multiple (MAE)-based approach for DeepJSCC design over arbitrary symmetric channels. again, a GAN is used in single eavesdropper. In this paper, we consider a AWGN winmap code based against multiple eavesdroppers. The multiple eavesdroppers for both colluding and non-colluding (MAE) eavesdroppers for AWGN JSCC with fading channels. For the scenario, of colluding eavesdroppers, the eavesdroppers share their logits to collaboratively infer private attributes based on the considered E2E training-based, while

\*Equal contribution









Hence, the following strategy is initially proposed for the system: The encoder and decoder functions of Alice and Bob should jointly minimize the S-LF, denoted by  $\mathcal{L}_{AB}$ . If often unknown, we estimate the expected distortion measure using samples  $\mathbf{u}_j$  from an available dataset  $\mathcal{D}_u$  by computing  $\mathbb{E}_{p(\mathbf{u}, \hat{\mathbf{u}})}[d(\mathbf{u}, f_{\Omega_S}(\mathbf{y}))] \approx \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} d(\mathbf{u}, \hat{\mathbf{u}})$ , where  $N_u = |\mathcal{D}_u|$ . 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 [1], [2] and wiretap channel [3], [4] employing adversarial likelihood compensation (ALC), which has been shown in [2] 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 minimize the systematics of adversarial predictions, resulting in the following loss function: 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  $\mathcal{L}_{AB}$ :

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  $\Delta^{\text{MSE}}(\mathbf{u}, \hat{\mathbf{u}})$ , and the structural similarity index (SSIM),  $\Delta^{\text{SSIM}}(\mathbf{u}, \hat{\mathbf{u}})$ , 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

The training process of legitimate nodes can be further enhanced via employing adversarial likelihood compensation (ALC), which has been shown in [2] 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:

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

where  $\mathcal{L}_{E,m}^{\text{MSE}}(\mathbf{u}, \hat{\mathbf{u}}) = \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} \mathbb{H}(q_{\Theta_{E,m}}(s|\mathbf{u}) | z_m(\mathbf{u}) | \mathbf{e}_s(\mathbf{u})) \triangleq \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} \mathbb{H}(q_{\Theta_{E,m}}(s|\mathbf{u}) | z_m(\mathbf{u}) | \mathbf{e}_s(\mathbf{u}))$ , and  $\alpha$  is a tuning parameter representing the contribution of the SSIM metric. Note that the adversarial SSIM can be trained in parallel. Then

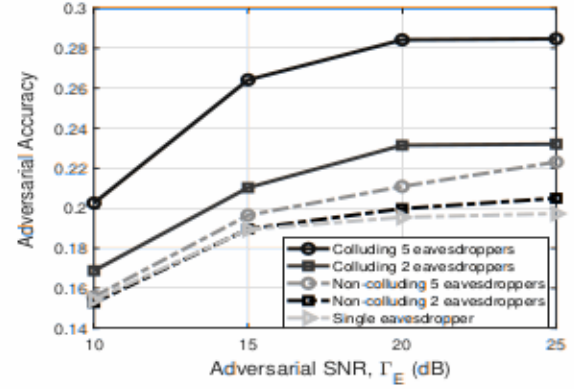


Fig. 4: Total adversarial accuracy over Rayleigh channels.

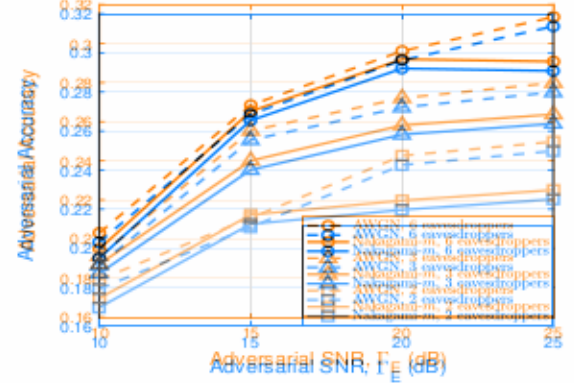


Fig. 5: Colluding adversarial accuracy over AWGN and Nakagami.

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, 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 [2]. 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 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.

Each step of the joint training of A-to-B autoencoder is followed by a training step for the adversarial DNNs. Eavesdroppers aim to minimize the CE between their estimated likelihood  $q_{\Theta_{E,m}}(s|\mathbf{u})$  and the ground-truth vector  $\mathbf{e}_s$  corresponding to  $S$ . Hence, the following LF is employed for training the DNN of Eve for  $m \in [M]$ :  $\mathcal{L}_{E,m}^{\text{MSE}}(\mathbf{u}, \hat{\mathbf{u}}) = \frac{1}{N_u} \sum_{\mathbf{u} \in \mathcal{D}_u} \mathbb{H}(q_{\Theta_{E,m}}(s|\mathbf{u}) | z_m(\mathbf{u}) | \mathbf{e}_s(\mathbf{u}))$ . Note that the adversarial DNNs can be trained by parallelizing the combination of the adversarial logits (the colluding step). The common secret  $S$  here is considered as the class of CIFAR-10 images with individually-extracted logits, and a weighted sum of these logits is exploited for the inference. We also set  $\alpha = 0.1$  and  $\beta = 1$  as default values. 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  $\Gamma_m = 10 \log_{10} \frac{P}{N} \text{ dB}$  and  $\Gamma_E = 10 \log_{10} \frac{P_E}{N_E} \text{ dB}$  (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 the proposed learning-based security solution. We also address the secrecy-utility trade-off for the proposed learning-based approach. For the training, we sample channel realizations from the general case of complex Rayleigh fading with average  $\Gamma_m$  and  $\Gamma_E$  SNR values. Nevertheless, during the inference phase, we study the performance in different scenarios of AWGN and Nakagami- $m$  channels (for  $m=2$ ). While we do assume known channel models in 1000 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 Python<sup>3</sup> 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 LF in the widely-adopted Adam optimizer, we choose [21] with a learning rate of  $10^{-4}$ . We fix the number of training episodes to  $N_{\text{episode}} = 200$  and the batch size to  $n_{\text{batch}} = 128$ . The DNNs with and without the adversarial accuracy, of the proposed scheme vs. the SNR of adversarial links (due to Eve) are shown in Fig. 3. In addition, the bandwidth compression ratio is set to 0.5 for the training, while the non-colluding benchmarks for 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 [21]. 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 (i.e.,  $s$ ). DNN architectures are implemented with the help of generative 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 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 LF in the widely-adopted Adam optimizer, we choose [21] with a learning rate of  $10^{-4}$ . We fix the number of training episodes to  $N_{\text{episode}} = 200$  and the batch size to  $n_{\text{batch}} = 128$ . The DNNs with and without the adversarial accuracy, of the proposed scheme vs. the SNR of adversarial links

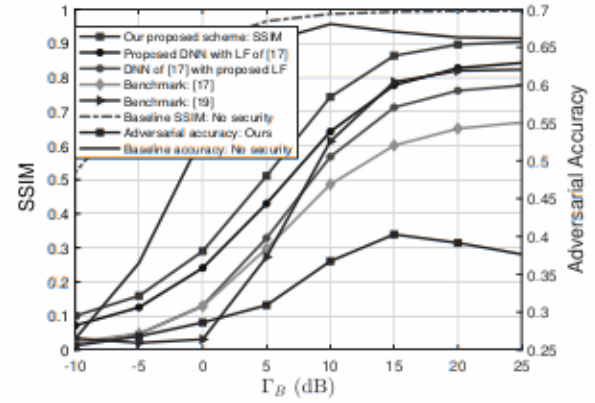


Fig. 3. Ablation study for  $M = 3$  non-colluding eavesdroppers.

(Fig. 3) illustrates the data reconstruction performance at Bob and the total adversarial accuracy, with having 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 [17] and [19], respectively. The ablation examination conducted in this figure show that both the implemented DNNs and the proposed LF for optimizing the framework contribute to the system's performance compared with other benchmarks. The figure also implies that increasing  $\Gamma_E$  results in having higher SSIM values. This is because increasing  $\Gamma_E$  can result in less distorted observations at Bob, which facilitates the image reconstruction performance. Data efficiency and adaptability of our proposed scheme are also validated, since we have trained our DNNs with a fixed SNR of 20 dB, while the performance gain of our approach during SNR inference holds for various SNR values. Fig. 3 highlights that if we ignore the eavesdroppers during the training of a pair and set the posterior adversarial distribution to  $\delta_{E,m}$ , the SSIM (i.e., data recovery) is impacted by the increase in  $\Gamma_E$  which highlights the limitation of this figure, where having 3 eavesdroppers can impact 10% decrease in the reconstruction performance of Bob. Finally, to indicate the importance of our proposed adversarial-aware scheme (in tem 3) 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 with part), the adversarial accuracy is increased by about 28%. This clearly highlights the importance of employing our proposed learning-based secure encoding scheme, when having  $M$  Fig. 3 illustrates the impact of parameter  $\gamma$  and the effect of adversarial purposes of our proposed system. These hyper-parameters of the coefficient  $\gamma$  are related with the training loss, adjusting  $\gamma$  can adjust the performance gain training loss by (23) and (22) respectively. In this experiment, the adversarial performance is investigated by the adversarial accuracy (i.e., the ratio of the ground truth and the predicted sensitive information) among the labels of CIFAR-10 dataset. The figure indicates that by increasing  $\gamma$ , higher values of SSIM are achieved, increasing  $\gamma$  emphasizes

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





- [48] M. P. Khatami and L. De Broeck, "Adaptive secret key generation with the-middle adversaries," *IEEE Wireless Comm. Lett.*, vol. 11, no. 4, pp. 856–860, Apr. 2022.
- [9] M. Letafati, H. Behroozi, B. H. Khalaj, and E. A. Jorswieck, "Deep learning for hardware-impaired wireless secret key generation with man-in-the-middle attacks," 2021 IEEE Global Communications Conference (GLOBECOM), Madrid, Spain, Dec. 2021, pp. 1–6.
- [10] E. Bourtsoulatz, D. Burth Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Trans. Cogn. Comm. Netw.*, vol. 5, no. 3, pp. 567–579, Sep. 2019.
- [11] D. B. Kurka and D. Gündüz, "Bandwidth-agile image transmission with deep joint source-channel coding," *IEEE Trans. Wireless Comm.*, vol. 20, no. 12, pp. 8081–8095, Dec. 2021.
- [12] M. Letafati, H. Behroozi, B. H. Khalaj, and E. A. Jorswieck, "Wireless-powered cooperative key generation for e-health: A reservoir learning approach," 2022 IEEE 95th Vehicular Technology Conference (VTC-Spring), Helsinki, Finland, Jun. 2022, pp. 1–7.
- [13] K. -L. Besser, P. -H. Lin, C. R. Janda, and E. A. Jorswieck, "Wiretap code design by neural network autoencoders," *IEEE Trans. Inf. Forensics Secur.*, vol. 15, pp. 3374–3386, Oct. 2019.
- [14] R. Fritschek, R. F. Schaefer, and G. Wunder, "Deep learning for channel coding via neural mutual information estimation," in *Proc. 20th Int. Workshop Sig. Proc. Adv. Wireless Comm.*, pp. 1–5, July 2019.
- [15] R. Fritschek, R. F. Schaefer, and G. Wunder, "Deep learning based wiretap coding via mutual information estimation," in *Proceedings of the 2nd ACM Workshop on Wireless Security and Machine Learning (WiseML'20)*, NY, USA, Jul. 2020, pp. 1–5.
- [16] T. Marchioro, N. Laurenti, and D. Gündüz, "Adversarial networks for secure wireless communications," *IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Barcelona, Spain, 2020, pp. 8748–8752.
- [17] E. Erdemir, P. L. Dragotti, and D. Gündüz, "Privacy-aware communication over a wiretap channel with generative networks," *IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Marina Bay Sands, Singapore, May 2022, pp. 2989–2993.
- [18] D. B. Kurka and D. Gündüz, "DeepJSCC-f: Deep joint source-channel coding of images with feedback," *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, pp. 178–193, May 2020.
- [19] R. Polikar, "Bootstrap-inspired techniques in computational intelligence," *IEEE Signal Process. Mag.*, vol. 24, no. 4, pp. 59–72, 2007.
- [20] A. Makhdoumi, S. Salamatian, N. Fawaz, and M. Médard, "From the information bottleneck to the privacy funnel," in *Proc. of IEEE Information Theory Workshop (ITW)*, 2014.
- [21] D. Barber and F. Agakov, "The IM algorithm: A variational approach to information maximization," *Advances in Neural Information Processing Systems 16: Proceedings of the 2003 Conference*, vol. 16, pp. 201–208.
- [22] H. Zhao, O. Gallo, I. Frosio, and J. Kautz, "Loss functions for image restoration with neural networks," *IEEE Transactions on Computational Imaging*, vol. 3, pp. 47–57, Mar. 2017.
- [23] A. Krizhevsky, "Learning multiple layers of features from tiny images," *University of Toronto, Tech. Rep.*, 2009.
- [24] D. P. Kingma and L. J. Ba, "Adam: A method for stochastic optimization," 2015 International Conference on Learning Representations (ICLR), San Diego, May 2015, pp. 1–13.