

WITT: A WIRELESS IMAGE TRANSMISSION TRANSFORMER FOR SEMANTIC COMMUNICATIONS

Ke Yang^{*}, Sixian Wang^{*}, Jincheng Dai^{*}, Kailin Tan^{*}, Kai Niu^{*†}, Ping Zhang^{*}

^{*} Beijing University of Posts and Telecommunications, Beijing, China

[†] Peng Cheng Laboratory, Shenzhen, China

Email: daijincheng@bupt.edu.cn

ABSTRACT

In this paper, we aim to redesign the vision Transformer (ViT) as a new backbone to realize semantic image transmission, termed wireless image transmission transformer (WITT). Previous works build upon convolutional neural networks (CNNs), which are inefficient in capturing global dependencies, resulting in degraded end-to-end transmission performance especially for high-resolution images. To tackle this, the proposed WITT employs Swin Transformers as a more capable backbone to extract long-range information. Different from ViTs in image classification tasks, WITT is highly optimized for image transmission while considering the effect of the wireless channel. Specifically, we propose a spatial modulation module to scale the latent representations according to channel state information, which enhances the ability of a single model to deal with various channel conditions. As a result, extensive experiments verify that our WITT attains better performance for different image resolutions, distortion metrics, and channel conditions. The code is available at <https://github.com/KeYang8/WITT>.

Index Terms— wireless image transmission, vision Transformer, joint source and channel coding

1. INTRODUCTION

Recently, learning-based joint source-channel coding (JSCC) for wireless data transmission emerges as an active research area in communication community [1–7]. By replacing the hand-crafted codecs with deep neural networks (DNNs), they achieve comparable or even better end-to-end transmission performance than traditional separation-based source and channel coding schemes. In particular, for image transmission tasks, deep JSCC [3] and its variants [4, 8–10] have competitive performance and much lower complexity compared to advanced image codec (JPEG/JPEG2000/BPG) followed by capacity-approaching channel code family (such as low-density parity-check (LDPC) coding [11]). Moreover, they can be agilely optimized for human visual perception [9], or downstream machine tasks [10]. Therefore, it is promising for many latency-sensitive applications, such as XR and autonomous driving.

Despite its great potential, previous works mainly build upon CNNs [3, 4, 8]. Limited by the model capacity, it can be observed that with the increase of the image dimension, the performance of the

CNN-based deep JSCC degrades rapidly and falls behind separation-based schemes. In this paper, we aim to break the aforementioned limits and increase the representation capacity of deep JSCC models. To this end, we expect to introduce the global attention mechanism among all the image patches to extract *high-level semantic features* of the source image for boosting a more efficient wireless image transmission method. Inspired by the recent advances of vision Transformer [12] in the computer vision field, it is the very time to redesign the vision Transformer as a new backbone for wireless image transmission.

In this paper, we propose a new JSCC framework named WITT, a high-efficiency wireless image transmission scheme that injects the advantages of the vision Transformer into the deep JSCC framework. By integrating the Swin Transformer [13] backbone in our scheme, a considerable performance gain can be achieved, especially for high-resolution images. Swin Transformer constructs hierarchical feature maps in the latent semantic space and has linear computational complexity to image size. Nevertheless, a naive change of the network backbone cannot obtain the expectable transmission performance gain over the imperfect wireless channels. To tackle this, we design a plug-in “Channel ModNet” inserted into Transformer to track the varying channel states. By this means, a single model can adapt to various channel states without retraining, which makes sense for the practical use of WITT.

We verify the performance of the proposed method through extensive experiments. We show that for image transmission, the proposed WITT method can achieve significant performance on various metrics such as PSNR and MS-SSIM [14]. Equivalently, the proposed method can save bandwidth costs by achieving identical end-to-end transmission performance. As the source image resolution increases, the performance superiority shows more clearly.

2. THE PROPOSED WITT SCHEME

2.1. Overall Architecture

An overview of the WITT architecture for wireless image transmission is given in Fig. 1(a). An RGB image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ is first split into $l_1 = \frac{H}{2} \times \frac{W}{2}$ non-overlapping patches. Each patch can be viewed as a “token”. We thus have a sequence of tokens (x_1, \dots, x_{l_1}) by putting these tokens in the order from top left to bottom right. After patch embedding, N_1 Swin Transformer blocks are applied on these l_1 tokens [13]. Here, we refer to these N_1 blocks together with the patch embedding layer as “stage 1”.

Then, these tokens are fed to several stages where “stage i ” is en-

This work was supported in part by the National Natural Science Foundation of China under Grant 92067202, Grant 62001049, Grant 62071058, and Grant 61971062, Beijing Natural Science Foundation under Grant 4222012, and the Fundamental Research Funds for the Central Universities.

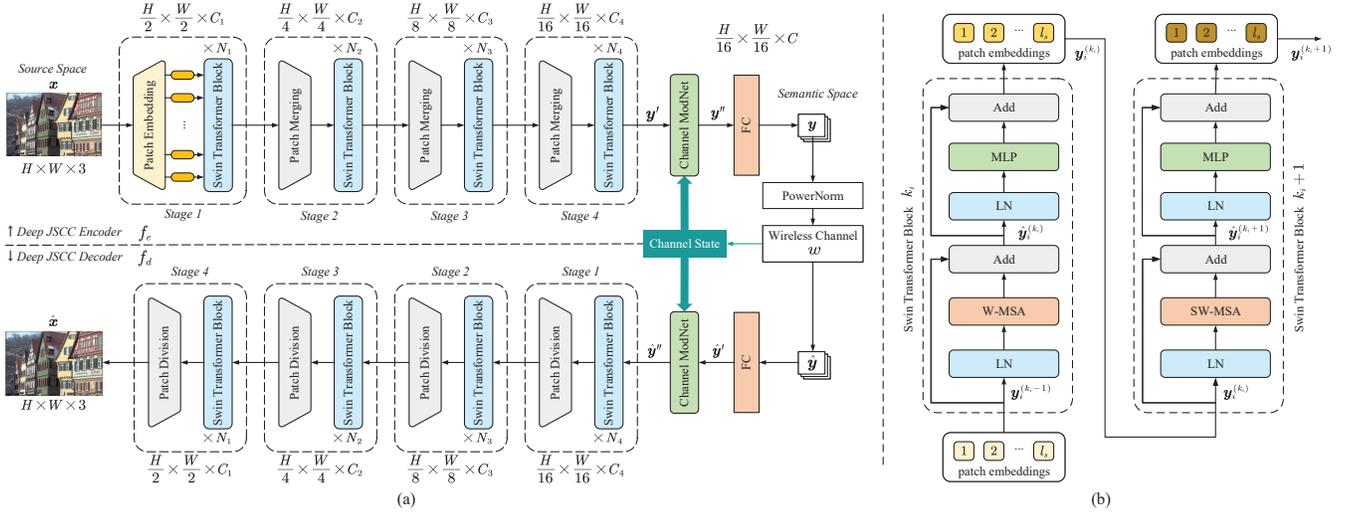


Fig. 1. (a) The overall architecture of the proposed WITT scheme for wireless image transmission. (b) Two successive Swin Transformer Blocks. W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

capsulated by a down-sampling patch merging layer and the following N_i Swin Transformer blocks. It is worth noting that the number of stages of the encoder f_e should vary with the image resolution. In general, images of higher resolution need more stages. As an example, Fig. 1(a) presents the four-stages version. After that, each patch embedding will be rescaled by a Channel ModNet according to the channel state. Next, an FC layer is applied on these embeddings to project it to C dimension. The channel bandwidth ratio is defined as $R = C / (2 \times 3 \times 2^n \times 2^n)$, where n denotes the number of stage.

Before transmitting \mathbf{y} into the wireless channel, the power normalization operation enables \mathbf{y} to satisfy the average power constraint. Then, the analog feature map is directly sent over the wireless channel. In this paper, we consider the general fading channel model with transfer function $\hat{\mathbf{y}} = w(\mathbf{y}; \mathbf{h}) = \mathbf{h} \odot \mathbf{y} + \mathbf{n}$, where \odot is the element-wise product, \mathbf{h} denotes the channel state information (CSI) vector, and each component of the noise vector \mathbf{n} is independently sampled from a Gaussian distribution, i.e., $\mathbf{n} \sim \mathcal{N}(0, \sigma_n^2 \mathbf{I}_k)$, where σ_n^2 is the average noise power.

The deep JSCC decoder f_d has a symmetric architecture with encoder f_e . It consists of an FC layer, Channel ModNet, patch division layers for up-sampling, and Swin Transformer blocks. It reconstructs input images from noisy latent representations $\hat{\mathbf{y}}$. The training loss function of the whole system is

$$\min_{\phi, \theta} \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \mathbb{E}_{\hat{\mathbf{y}} \sim p_{\hat{\mathbf{y}}|\mathbf{x}}} [d(\mathbf{x}, \hat{\mathbf{x}})], \quad (1)$$

where $\hat{\mathbf{y}} = w(f_e(\mathbf{x}; \phi); \nu)$, $\hat{\mathbf{x}} = f_d(\hat{\mathbf{y}}; \theta)$, ϕ and θ encapsulate all the network parameters of f_e and f_d , respectively.

2.2. Swin Transformer Block

As shown in Fig. 1(b), a Swin Transformer block is a sequence-to-sequence function that is built by replacing the standard multi-head self attention (MSA) module in a Transformer block with a module based on shifted windows [13]. The shift of the window partition between consecutive self attention layers provides connections among them, significantly enhancing modeling power.

With the shifted window partitioning approach, consecutive Swin Transformer blocks k_i and $k_i + 1$ of “stage i ” are computed as

$$\hat{\mathbf{y}}_i^{(k_i)} = \text{W-MSA}(\text{LN}(\mathbf{y}_i^{(k_i-1)})) + \mathbf{y}_i^{(k_i-1)}, \quad (2a)$$

$$\mathbf{y}_i^{(k_i)} = \text{MLP}(\text{LN}(\hat{\mathbf{y}}_i^{(k_i)})) + \hat{\mathbf{y}}_i^{(k_i)}, \quad (2b)$$

$$\hat{\mathbf{y}}_i^{(k_i+1)} = \text{SW-MSA}(\text{LN}(\mathbf{y}_i^{(k_i+1)})) + \mathbf{y}_i^{(k_i+1)}, \quad (3a)$$

$$\mathbf{y}_i^{(k_i+1)} = \text{MLP}(\text{LN}(\hat{\mathbf{y}}_i^{(k_i+1)})) + \hat{\mathbf{y}}_i^{(k_i+1)}, \quad (3b)$$

where $\hat{\mathbf{y}}_i^{(k_i)}$ and $\mathbf{y}_i^{(k_i)}$ represent the output feature of the (S)W-MSA module and the MLP module for block k_i at stage i , respectively. W-MSA and SW-MSA are multi-head self attention modules with regular and shifted operation configurations. LN denotes the layer normalization operation [12].

2.3. Channel ModNet

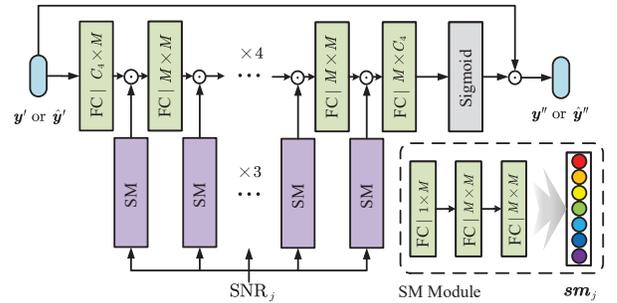


Fig. 2. The architecture of Channel ModNet. C_4 and M denote the number of channels in \mathbf{y}' or $\hat{\mathbf{y}}$ and the number of intermediates of FCN in \mathbf{sm}_j respectively.

The proposed “Channel ModNet” is a plug-in module to modulate the output of several Transformer stages as shown in Fig. 1(a). For different channel states, Channel ModNet can generate specific deep JSCC codec functions to adapt to channel changes.

In particular, the semantic feature map \mathbf{y}' is fed into our Channel ModNet to be modulated by the channel state information. Correspondingly, the received symbol $\hat{\mathbf{y}}$ is first processed by the FC layer and then sent into our ModNet. Thus, for both f_e and f_d , the channel state is taken as a coding factor sent into Channel ModNet,

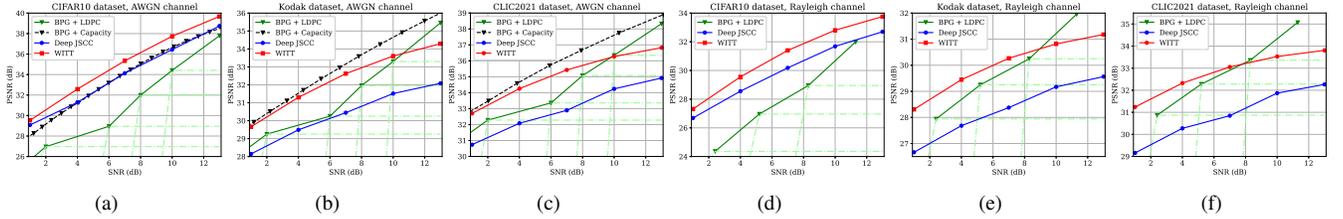


Fig. 3. (a)~(c) PSNR performance versus the SNR over the AWGN channel. (d)~(f) PSNR performance versus the SNR over the Rayleigh fast fading channel. The average CBR is set to 1/3, 1/16, and 1/16 for CIFAR10 dataset, Kodak dataset, and CLIC2021 dataset.

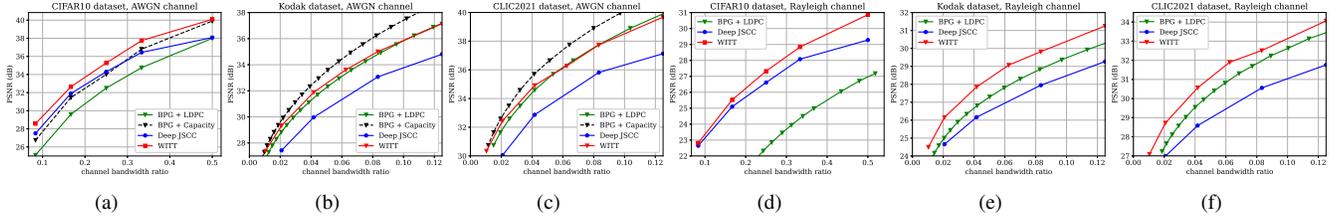


Fig. 4. (a)~(c) PSNR performance versus the CBR over the AWGN channel at SNR = 10dB. (d)~(f) PSNR performance versus the CBR over the Rayleigh fast fading channel at SNR = 3dB.

which modulates the intermediate feature maps to the wireless channel state.

The architecture of the Channel ModNet is depicted in Fig. 2. It consists of 8 FC layers alternating with 7 SNR modulation (SM) modules. SM module is a three-layered FC network, which transforms the input SNR_j into an M -dimensional vector \mathbf{sm}_j . Multiple SM modules are cascaded sequentially in a coarse-to-fine manner. The previous modulated features are fed into subsequent SM modules. The arbitrary target modulator can be realized by assigning a corresponding SNR value. The mapping procedures from SNR_j to \mathbf{sm}_j are

$$\mathbf{sm}_j^{(1)} = \text{ReLU}(\mathbf{W}^{(1)} \cdot \text{SNR}_j + \mathbf{b}^{(1)}), \quad (4a)$$

$$\mathbf{sm}_j^{(2)} = \text{ReLU}(\mathbf{W}^{(2)} \cdot \mathbf{sm}_j^{(1)} + \mathbf{b}^{(2)}), \quad (4b)$$

$$\mathbf{sm}_j = \text{ReLU}(\mathbf{W}^{(3)} \cdot \mathbf{sm}_j^{(2)} + \mathbf{b}^{(3)}), \quad (4c)$$

where ReLU and Sigmoid are the activation functions, \mathbf{W} and \mathbf{b} are the affine function parameters and their corresponding bias.

Therefore, the coding SNR_j is associated with a tensor \mathbf{sm}_j in each SM module. Then, the input feature will be fused with \mathbf{sm}_j in the element-wise product, i.e.,

$$\mathbf{output} = \mathbf{input} \odot \mathbf{sm}_j \quad (5)$$

Here, \mathbf{input} denotes the feature output from the previous FC layer, and \mathbf{output} is feeding into the next FC layer.

3. EXPERIMENTAL RESULTS

3.1. Experimental Setup

Datasets: We train and evaluate the proposed WITT scheme on image datasets with different resolutions from 32x32 upto 2K. For low-resolution images, we use the CIFAR10 [15] dataset for training and testing. For high-resolution images, we choose DIV2K [16] dataset for training, and use the Kodak [17] dataset and the CLIC2021 [18] testset for testing. During training, images are randomly cropped into 256×256 patches.

Comparison Schemes: We compare our WITT scheme with the

CNN-based deep JSCC scheme [3] and classical separation-based source and channel coding schemes. Specifically, we employ the BPG [19] codec for compression combined with 5G LDPC codes [11] for channel coding (marked as “BPG + LDPC”). Here, we considered 5G LDPC codes with a block length of 6144 bits for different coding rates and quadrature amplitude modulations (QAM). Moreover, the ideal capacity-achieving channel code is also considered during evaluation (marked as “BPG + Capacity”).

Evaluation Metrics: We qualify the performance of the proposed scheme using both the widely used pixel-wise metric PSNR and the perceptual metric MS-SSIM [20]. For PSNR, we optimized our model by the mean square error (MSE) loss function. For MS-SSIM, the loss function is $1 - \text{MS-SSIM}$.

Training Details: The number of stages in WITT varies with training image resolution. For low-resolution images, we use 2 stages with $[N_1, N_2] = [2, 4]$, $[C_1, C_2] = [128, 256]$, and the window size is set to 2. For large-resolution images, we use 4 stages $[N_1, N_2, N_3, N_4] = [1, 1, 2, 6]$, $[C_1, C_2, C_3, C_4] = [128, 192, 256, 320]$, and the window size is set to 8. During the training process, we first train other parameters except for the Channel ModNet over the wireless channel. Then, the whole proposed model is trained with Channel ModNet. We exploit the Adam optimizer with a learning rate of 1×10^{-4} , and the batch size is set to 128 and 16 for CIFAR10 dataset and DIV2K dataset, respectively. The WITT model is trained under the channel with a uniform distribution of $\text{SNR}_{\text{train}}$ from 1dB to 13 dB.

3.2. Result Analysis

Fig. 3(a) ~ 3(c) show the PSNR performance versus SNR over the AWGN channel, and Fig. 3(d) ~ 3(f) present the Rayleigh fast fading channel. For the WITT scheme, a single model can cover a range of SNR from 1dB to 13dB. For the “BPG + LDPC” scheme, according to adaptive modulation and coding (AMC) standard [21], we choose the best-performing configuration of coding rate and modulation (the green dashed lines) under each specific SNR and plot the envelope. Compared to the CNN-based deep JSCC

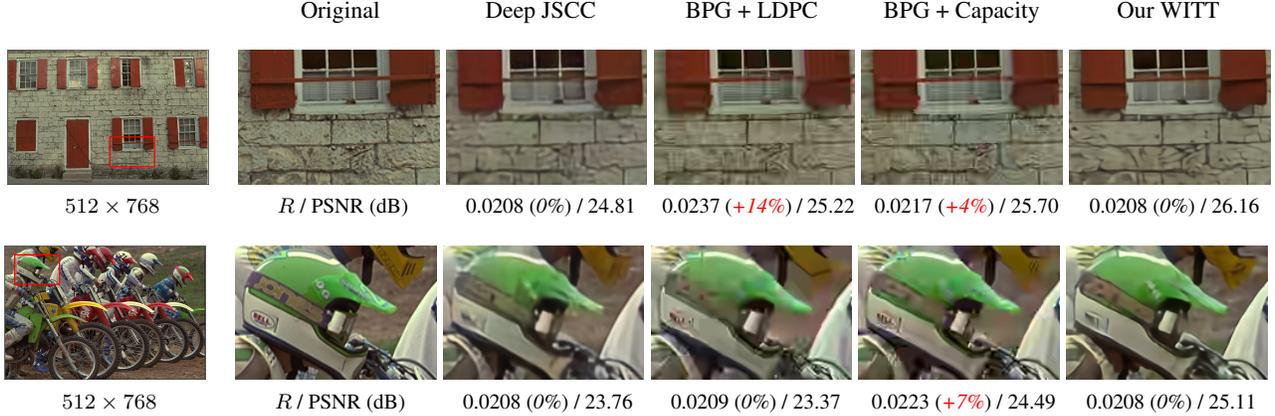


Fig. 5. Examples of visual comparison under AWGN channel at SNR = 10dB. The first column, second column, and third to sixth column shows the original image, original patch, and reconstructions of different transmission schemes, respectively. The red number indicates the percentage of extra bandwidth cost compared to WITT.

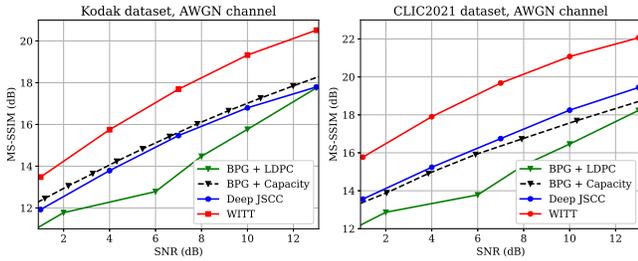


Fig. 6. MS-SSIM performance versus the SNR at the AWGN channel, and the average CBR is set to 1/16.

scheme, we achieve much better performance for all SNRs. Due to the enhanced model capacity by incorporating Transformers, it can be seen that the performance gap increases with the growth of image resolution. For the CIFAR10 dataset, WITT and deep JSCC scheme significantly outperform the “BPG + LDPC” and “BPG + Capacity”. However, for high-resolution images, the performance of CNN-based deep JSCC degrades a lot and falls behind to the separation-based scheme. Our proposed maintains a considerable performance, especially in the low SNR regions.

Fig. 4(a) ~ 4(c) demonstrate the PSNR performance versus the CBR over the AWGN channel, and Fig. 4(d) ~ 4(f) show the Rayleigh fast fading channel. For the CIFAR10 dataset, our proposed model can generally outperform deep JSCC for all CBRs. Meanwhile, our model achieves considerable gains compared to the existing classical separation-based schemes, especially on the Rayleigh channel. For high-resolution image datasets, our proposed model outperforms the CNN-based deep JSCC scheme. Compared to “BPG + LDPC”, our WITT model achieves comparable or better performance and coding gain. Moreover, WITT cannot provide comparable coding gain as that of the “BPG + Capacity” scheme, i.e., the slope of the performance curve slows down with the increase of SNR. Nevertheless, the performance of WITT approaches the “BPG + Capacity” in the low CBR regions and obviously improves compared to CNN-based deep JSCC.

Besides, to more comprehensively evaluate the performance of our model, we also train and test our model on the MS-SSIM metric.

Table 1. Inference speed and complexity comparison.

method	inference time	FLOPs	#param.
WITT	116ms	198G	28.2M
ADJSCC	155ms	511G	16.2M

MS-SSIM is a multi-scale perceptual metric that approximates human visual perception well. Fig. 6 shows the performance versus the SNR at the AWGN channel with an average CBR = 1/16. For more intuitive observation and comparison, it is converted into the form of dB and the formula is $MS-SSIM(dB) = -10\log(1 - MS-SSIM)$. Results indicate that the proposed WITT model can outperform other competitors by a large margin. Compared to the PSNR results in Fig. 3, we can find that classical image transmission series are inferior to the learning-based WITT because classical image compression is designed to be optimized for squared error with hand-crafted constraints. Fig. 5 visualizes the reconstructions. It can be observed that WITT can achieve better visual quality with the same or lower channel bandwidth cost. More specifically, it avoids block artifacts and produces higher fidelity textures and details.

Table 1 lists the inference time, FLOPs, and model size (#param.) for WITT and ADJSCC [4] on Kodak dataset with a batch size of 1. All experiments are carried out using a Linux server with a single RTX 3090 GPU. Benefiting from the high-efficient window-based attention mechanism, WITT spends 2.5x lower floating point of operations (FLOPs). Despite its larger model size, WITT can provide better performance and run faster than the ADJSCC.

4. CONCLUSION

In this paper, we have proposed a high-efficiency scheme named WITT to improve the performance of wireless image transmission. The WITT framework is built upon the Swin Transformer to extract long-term hierarchical image representation. To deal with various channel conditions, we have further proposed the Channel ModNet to rescale the representations according to channel states automatically. Results have demonstrated that our proposed method outperforms the CNN-based deep JSCC scheme and the classical separated-based schemes.

5. REFERENCES

- [1] K. Choi, K. Tatwawadi, A. Grover, T. Weissman, and S. Ermon, "Neural joint source-channel coding," in *International Conference on Machine Learning*. PMLR, 2019, pp. 1182–1192.
- [2] N. Farsad, M. Rao, and A. Goldsmith, "Deep learning for joint source-channel coding of text," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2018, pp. 2326–2330.
- [3] E. Bourtsoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, 2019.
- [4] J. Xu, B. Ai, W. Chen, A. Yang, P. Sun, and M. Rodrigues, "Wireless image transmission using deep source channel coding with attention modules," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [5] H. Xie, Z. Qin, G. Y. Li, and B. H. Juang, "Deep learning enabled semantic communication systems," *IEEE Transactions on Signal Processing*, vol. 69, pp. 2663–2675, 2021.
- [6] J. Dai, P. Zhang, K. Niu, S. Wang, Z. Si, and X. Qin, "Communication beyond transmitting bits: Semantics-guided source and channel coding," *IEEE Wireless Communications*, pp. 1–8, 2022.
- [7] P. Zhang, W. Xu, H. Gao, K. Niu, X. Xu, X. Qin, C. Yuan, Z. Qin, H. Zhao, J. Wei, et al., "Toward wisdom-evolutionary and primitive-concise 6g: A new paradigm of semantic communication networks," *Engineering*, vol. 8, pp. 60–73, 2022.
- [8] 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, 2020.
- [9] J. Wang, S. Wang, J. Dai, Z. Si, D. Zhou, and K. Niu, "Perceptual learned source-channel coding for high-fidelity image semantic transmission," *Proceedings of IEEE Global Communications Conference*, 2022.
- [10] M. Jankowski, D. Gündüz, and K. Mikołajczyk, "Wireless image retrieval at the edge," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 89–100, 2020.
- [11] T. Richardson and S. Kudekar, "Design of low-density parity-check codes for 5G new radio," *IEEE Communications Magazine*, vol. 56, no. 3, pp. 28–34, 2018.
- [12] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [13] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9992–10002.
- [14] K. Ding, K. Ma, S. Wang, and E. P. Simoncelli, "Comparison of full-reference image quality models for optimization of image processing systems," *International Journal of Computer Vision*, vol. 129, no. 4, pp. 1258–1281, 2021.
- [15] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Master's thesis, University of Tront*, 2009.
- [16] E. Agustsson and R. Timofte, "Ntire 2017 challenge on single image super-resolution: Dataset and study," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 1122–1131.
- [17] "Kodak PhotoCD dataset," URL: <http://r0k.us/graphics/kodak/>, 1993.
- [18] "CLIC 2021: Challenge on learned image compression," URL: <http://compression.cc>, 2021.
- [19] F. Bellard, "BPG image format," URL: <https://bellard.org/bpg/>.
- [20] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers*. IEEE, 2003, vol. 2, pp. 1398–1402.
- [21] 3GPP, "NR; Physical layer procedures for data," Technical Specification (TS) 38.214, 3rd Generation Partnership Project (3GPP), 2018, Version 15.0.0.