

# Semantic Communications with Explicit Semantic Base for Image Transmission

Yuan Zheng\*, Fengyu Wang<sup>†</sup>, Wenjun Xu\*, Miao Pan<sup>‡</sup>, and Ping Zhang\*

\*State Key Laboratory of Network and Switching Technology,

Beijing University of Posts and Telecommunications, Beijing, 100876, China

<sup>†</sup>School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, 100876, China

<sup>‡</sup>Department of Electrical and Computer Engineering, University of Houston, Houston, TX 77004, USA

*Abstract*—Semantic communications, aiming at ensuring the successful delivery of the meaning of information, are expected to be one of the potential techniques for the next generation communications. However, the knowledge forming and synchronizing mechanism that enables semantic communication systems to extract and interpret the semantics of information according to the communication intents is still immature. In this paper, we propose a semantic image transmission framework with explicit semantic base (Seb), where Sebs are generated and employed as the knowledge shared between the transmitter and the receiver with flexible granularity. To represent images with Sebs, a novel Seb-based reference image generator is proposed to generate Sebs and then decompose the transmitted images. To further encode/decode the residual information for precise image reconstruction, a Seb-based image encoder/decoder is proposed. The key components of the proposed framework are optimized jointly by end-to-end (E2E) training, where the loss function is dedicatedly designed to tackle the problem of non-differentiable operation in Seb-based reference image generator by introducing a gradient approximation mechanism. Extensive experiments show that the proposed framework outperforms state-of-art works by 0.5 – 1.5 dB in peak signal-to-noise ratio (PSNR) w.r.t. different signal-to-noise ratios (SNR).

*Index Terms*—semantic communication system, semantic base, image transmission

## I. INTRODUCTION

Over the past decades, wireless communications have developed rapidly, where the global mobile traffic is expected to explode up to thousands of exabytes (EB) per month by 2030 [1]. However, the conventional communication systems based on the classic information theory focus on the bit-level precise transmission while ignoring the meaning of information. With channel capacity and source coding efficiency approaching the Shannon’s limit [2], [3], the conventional communications cannot meet the upsurging communication demands brought by the intelligent services such as Industrial Internet of Things (IIoT), Internet of Vehicles (IoV), and Extended Reality (XR).

Recently, benefiting from the great success of deep learning (DL), semantic communications, aiming to successfully convey the meaning of information between transceivers, have experienced vigorous development. Most existing works develop a synchronized knowledge base between the transmitter

and the receiver to support the extraction and interpretation of semantics according to the communication intents [4]–[7], in which knowledge sharing is indispensable. However, most of the current works [8]–[11] take the parameters of the proposed neural network as background knowledge, where the knowledge sharing is implicitly applied in the process of encoding and decoding. Note that the implicit knowledge sharing scheme lacks the awareness of basic features of semantic representation (e.g. the granularity and efficiency), which can significantly affect the performance of the entire system. To overcome the problem, the concept of semantic base (Seb) is proposed [12], which can be delicately designed with different levels of granularity and similarity to represent semantics w.r.t. intents of communications. The potential of semantic communications is expected to be further unleashed by properly generating Sebs to form the shared knowledge base between transceivers. However, effective Seb generating schemes have not been developed yet.

To fill the gap, in this paper, a Seb-based semantic communication framework for image transmission is proposed. Taking the complex correlated image sets as input, the proposed framework first splits the image sets according to the correlation level, constructs subsets containing more stably related images for further processing, and generates Sebs to represent the shared semantics as synchronized knowledge. During the transmission, each image is represented by the generated Sebs, and the residual information is further encoded for precise recovery of the image. The contributions of this paper are summarized as follows.

- A novel semantic communication framework with explicit Sebs for image transmission is proposed, where explicit procedures of Seb generation, transmission, and image semantics representation based on Sebs are included. To the best of our knowledge, this is the first work that utilizes explicit Sebs for image transmission.
- For Seb generation and Seb-based image representation, the **Seb-based reference image generator** is dedicatedly designed. The **Seb-based image encoder/decoder** is further proposed to encode/decode the residual information for precise image reconstruction. A specialized loss function with a gradient approximation mechanism is introduced to enable the E2E training for the entire network

This work is supported, in part, by the National Natural Science Foundation of China under Grant 62293485, and in part by the Fundamental Research Funds for the Central Universities under Grant 2022RC18. Corresponding author: Fengyu Wang, Email: fengyu.wang@bupt.edu.cn.

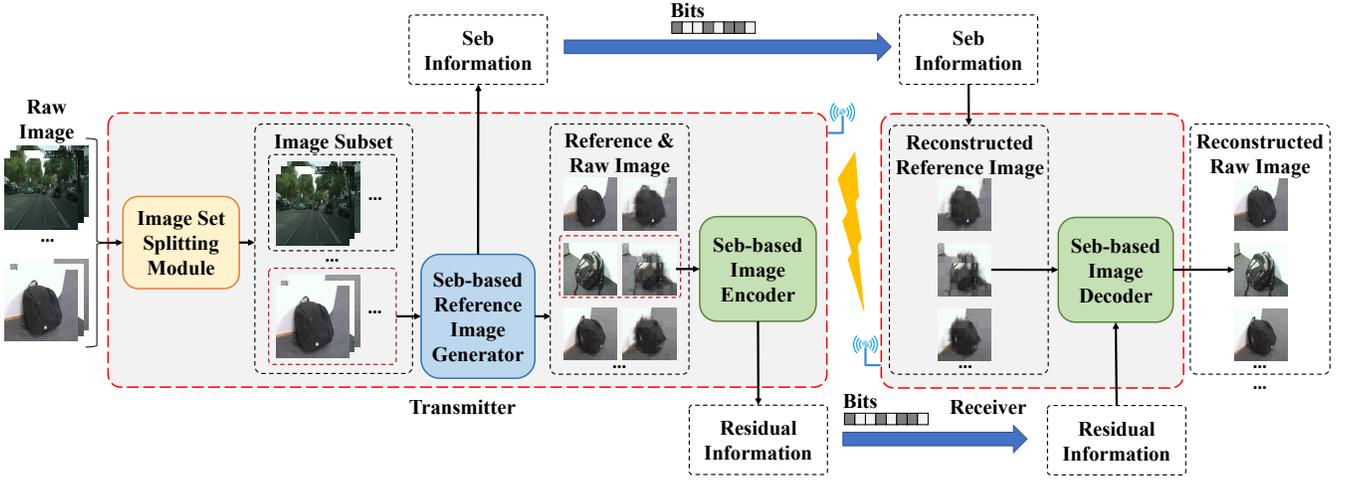


Fig. 1. The Seb-based image transmission framework.

that contains a non-derivable clustering algorithm in the **Seb-based reference image generator**.

- Extensive experiments validate the effectiveness and the generality of the proposed framework, which outperforms traditional methods and state-of-the-art DL-based methods by more than 0.5 – 1.5 dB w.r.t. peak signal-to-noise ratio (PSNR) under the same signal-to-noise ratio (SNR) and channel bandwidth ratio (CBR).

## II. SEB-BASED IMAGE TRANSMISSION FRAMEWORK

The Seb-based semantic communication framework is shown in Fig. 1. The transmitter includes an **image set splitting module**, a **Seb-Based reference image generator**, and a **Seb-based image encoder**. The corresponding Sebs and residual of raw images are extracted and then encoded by a channel encoder for transmission. At the receiver side, the **Seb-based image decoder** reconstructs raw images according to Sebs and their corresponding residual information. The detailed design of each module will be introduced in the following.

### A. Image Set Splitting Module

The image set splitting module is employed to split the original image set into several subsets, in which the images are expected to possess higher and more stable correlations that can be better captured by subsequent modules. Denoting  $\mathcal{I} = \{I_1, I_2, \dots, I_n\}$ ,  $I \in \mathbb{R}^{C \times H \times W}$  as the original image set, which is composed of  $n$  images  $I_1, I_2, \dots, I_n$ , where channels, height, and width of each image are denoted as  $C$ ,  $H$ , and  $W$ , respectively. We follow the architecture of contrastive learning [13] to build the image set splitting module, where

$$\{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_J\} = f_{\text{cluster}}(f_{\alpha}(\mathcal{I}, \alpha), J). \quad (1)$$

$f_{\alpha}(\cdot, \alpha)$  denotes the projector with its parameters as  $\alpha$  that projects each image into latent space, where similar images are projected into the same cluster. The projector is achieved by a Resnet-50 backbone and a projection multilayer perceptron

(MLP) in this paper.  $f_{\text{cluster}}(\cdot, J)$  denotes the standard k-means clustering algorithm with  $J$  as the number of clusters, which is determined by the clustering inertia. The module outputs  $J$  image subsets  $\mathcal{I}_j = \{I_{j(1)}, \dots, I_{j(n_j)}\}$ ,  $j = 1, \dots, J$  with each subset  $\mathcal{I}_j$  consisting of  $n_j$  images individually. The projector is initialized with the weights in [13], and is frozen as an invariant image classifier during training.

### B. Seb-based Reference Image Generator

To achieve Seb generation and representation of images, the Seb-based reference image generator is employed mainly including a Seb generator and a reference patch generator, as shown in Fig. 2.

In specific, the image subset  $\mathcal{I}_j$  is first divided into the patch set  $\mathcal{P}_j = \{P_{j(1)(1)}, \dots, P_{j(1)(n_p)}, \dots, P_{j(n_j)(n_p)}\}$ , where  $P \in \mathbb{R}^{C \times h \times w}$ . Each image is divided into  $n_p$  patches (e.g.,  $I_{j(1)}$  corresponds to  $\{P_{j(1)(1)}, \dots, P_{j(1)(n_p)}\}$ ). Parameters  $h$  and  $w$  denote the patch height and patch width, respectively. The operation brings about the reduction of size, and as a result, the patches will be mapped into a latent space with lower dimensions so that the semantics will be more precisely represented by Sebs. The granularity of Seb is thus controlled by the patch number and the patch size.

The Seb generator and the reference patch generator are designed based on the autoencoder structure, which are composed of convolutions and generalized divisive normalization (GDN)/inverse generalized divisive normalization (IGDN) [14] activation functions. The Seb generator is first used to extract semantic features  $\mathcal{F}_j = \{F_{j(1)(1)}, \dots, F_{j(1)(n_p)}, \dots, F_{j(n_j)(n_p)}\}$  from the patches  $\mathcal{P}_j$ , where  $F \in \mathbb{R}^{c' \times h' \times w'}$ , with  $c' = \frac{c}{16}$ , and  $w' = \frac{w}{16}$  denoting the dimensions of latent.

Next, the set of Sebs  $\mathcal{S}_j = \{S_{j(1)}, \dots, S_{j(K)}\}$ , where  $S \in \mathbb{R}^{c' \times h' \times w'}$ , and the corresponding usage information  $\mathcal{A}_j = \{A_{j(1)(1)}, \dots, A_{j(1)(n_p)}, \dots, A_{j(n_j)(n_p)}\}$ , where  $A \in \{1, \dots, K\}$  are generated by a standard clustering algorithm (e.g. K-means) on the semantic features  $\mathcal{F}_j$ . Each Seb  $S \in \mathcal{S}_j$  corresponds to the center of each cluster, where the usage

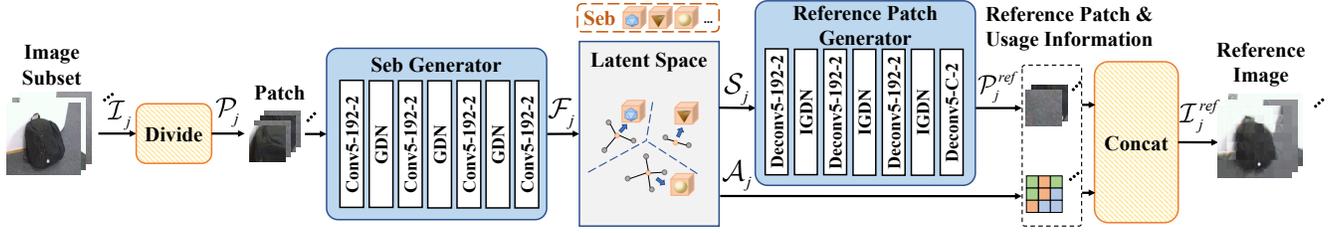


Fig. 2. The Seb-based reference image generator. Conv/Deconv5-192-2 represents the convolution/deconvolution operation with  $5 \times 5$  kernel size, 192 output channels, and a stride of 2. GDN/IGDN [14] denotes the nonlinear transform function. The number of Seb  $S_j$  is much smaller than the number of the latent feature of image patches  $F_j$  ( $K \ll n_j * n_p$ ), which supports efficient semantic representation. The whole structure is deployed at the transmitter.

information  $A \in \mathcal{A}_j$  records the index of the cluster to which the corresponding semantic feature  $F \in \mathcal{F}_j$  belongs. Parameter  $K \in \mathbb{Z}^+$  denotes the number of clustering centers, which controls the representation efficiency of Sebs with  $K \ll n_j * n_p$ .

To represent images with Sebs, the set of reference patches  $\mathcal{P}_j^{ref} = \{P_{j(1)}^{ref}, \dots, P_{j(K)}^{ref}\}$ , where  $P^{ref} \in \mathbb{R}^{C \times h \times w}$  are generated based on Sebs  $S_j$  by mapping the latent into the original space through the reference image generator. The reference image  $\mathcal{I}_j^{ref} = \{I_{j(1)}^{ref}, \dots, I_{j(n_j)}^{ref}\}$ ,  $I^{ref} \in \mathbb{R}^{C \times H \times W}$  is generated through a concatenating operation under the guidance of Seb usage information  $\mathcal{A}_j$  and the corresponding patches  $\mathcal{P}_j^{ref}$ . The mechanism of the **Seb-based reference image generator** is described as follows:

$$\begin{aligned} \mathcal{P}_j &= f_{\text{divide}}(\mathcal{I}_j), \\ (\mathcal{S}_j, \mathcal{A}_j) &= f_{\text{cluster}}(f_{\phi}(\mathcal{P}_j; \phi), K), \\ \mathcal{I}_j^{ref} &= f_{\text{concat}}(f_{\theta}(\mathcal{S}_j; \theta), \mathcal{A}_j), \end{aligned} \quad (2)$$

where  $f_{\text{divide}}(\cdot)$ ,  $f_{\text{cluster}}(\cdot)$ , and  $f_{\text{concat}}(\cdot)$  denote the fore-mentioned dividing, clustering, and concatenating operations, respectively,  $f_{\phi}(\cdot; \phi)$  and  $f_{\theta}(\cdot; \theta)$  denote the corresponding Seb generator and reference patch generator with  $\phi$  and  $\theta$  standing for their trainable parameters.

Note that Sebs  $S_j$  and the usage information  $\mathcal{A}_j$  need to be synchronized between the transmitter and the receiver. The wireless channel is modeled as

$$\begin{aligned} y &= hx + z, \\ \hat{x} &= h^{-1}y = x + \hat{z}, \end{aligned} \quad (3)$$

where  $x \in \mathbb{C}$  and  $y \in \mathbb{C}$  denote the complex symbols of channel input and output,  $z \sim \mathcal{CN}(0, \sigma^2)$  denotes the additive white Gaussian noise (AWGN) with  $\sigma^2$  as the average noise power, and  $h \in \mathbb{C}$  denotes the channel gain. At the receiver, by estimating the channel state information (CSI), the recovery of the channel input  $\hat{x} \in \mathbb{C}$  can be obtained, where  $\hat{z} = h^{-1}z$ .

$S_j$  and  $\mathcal{A}_j$  are transmitted after mapping into complex symbols with channel coding and modulation, and are recovered as  $\hat{S}_j$  and  $\hat{\mathcal{A}}_j$  at the receiver. The concatenation operation is carried out at the receiver as well to obtain  $\hat{\mathcal{I}}_j^{ref}$ , where

$$\hat{\mathcal{I}}_j^{ref} = f_{\text{concat}}(f_{\theta}(\hat{S}_j; \theta), \hat{\mathcal{A}}_j). \quad (4)$$

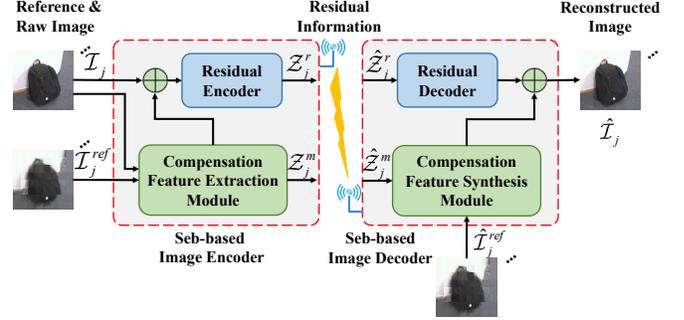


Fig. 3. The Seb-based image encoder/decoder.

### C. Seb-based Image Encoder/Decoder

To achieve precise image reconstruction, the **Seb-based image encoder/decoder** is proposed to extract residual information based on the reference images  $\mathcal{I}_j^{ref}$ . It includes a pair of compensation feature extraction/synthesis module and a pair of residual encoder/decoder, as shown in Fig. 3.<sup>1</sup>

In the encoding process, the compensation feature extraction module takes the raw images  $\mathcal{I}_j$  and the generated reference images  $\mathcal{I}_j^{ref}$  as inputs, and uses a CNN-based optical flow net [17] and a corresponding encoder to extract the compensation features  $\mathcal{Z}_j^m$  that corresponds to the predictions to the raw images. The residual encoder further compresses the residual parts between the raw images and their predictions into the residual features  $\mathcal{Z}_j^r$ . The features  $\mathcal{Z}_j^m$  and  $\mathcal{Z}_j^r$  need to be transmitted through the wireless channel as depicted in (3), and be recovered as  $\hat{\mathcal{Z}}_j^m$  and  $\hat{\mathcal{Z}}_j^r$  at the receiver. In the decoding process, the residual parts and the predictions of the images are obtained from the recovered  $\hat{\mathcal{Z}}_j^m$  and  $\hat{\mathcal{Z}}_j^r$  by the residual decoder and the compensation feature synthesis module, respectively. Then the reconstructed image set  $\hat{\mathcal{I}}_j$  is obtained through a summation operation.

## III. PROBLEM FORMULATION AND TRAINING STRATEGY

The key components of the proposed framework are trained in an E2E manner, and are jointly optimized based on rate-

<sup>1</sup>We follow the structure of E2E deep video compression scheme [16] to build the proposed Seb-based image encoder/decoder in this paper. The structure of the Seb-based image encoder/decoder will be further optimized in the future journal manuscript.

distortion trade-off through a specialized loss function.

### A. Rate-Distortion Optimization

The goal of the proposed Seb-based image transmission framework is to minimize the number of transmitted bits, while making the distortion between the raw image set  $\mathcal{I}$  and the reconstructed image set  $\hat{\mathcal{I}}$  as small as possible. This is a typical rate-distortion optimization problem that is modeled as

$$L_{RD} = \lambda D(\mathcal{I}, \hat{\mathcal{I}}) + (R(\mathcal{S}) + R(\mathcal{A}) + R(\mathcal{Z}^m) + R(\mathcal{Z}^r)). \quad (5)$$

$D(\mathcal{I}, \hat{\mathcal{I}})$  denotes the distortion between the raw and reconstructed image set, which is represented by the mean square error (MSE) in training.  $R(\mathcal{S}), R(\mathcal{A}), R(\mathcal{Z}^m), R(\mathcal{Z}^r)$  denote the bitrate corresponding to the Seb, the Seb usage information, as well as the compensation and the residual features, respectively.  $\lambda$  is a hyper-parameter that controls the trade-off between rate and distortion. Higher  $\lambda$  frameworks tend to consume more resources for less distortion.

In specific, the bitrate of Seb for the image set  $\mathcal{I}$  is approximated as

$$R(\mathcal{S}) = \sum_{S_j \in \mathcal{S}} R(S_j) = \sum_{S_j \in \mathcal{S}} \sum_{Pref \in \mathcal{P}^{ref}} R(P^{ref}), \quad (6)$$

where  $R(P^{ref})$  denotes the transmitting cost of  $P^{ref}$ . The bitrate of Seb usage information  $\mathcal{A}$  is calculated as

$$R(\mathcal{A}) = \sum_{A_j \in \mathcal{A}} R(A_j) = \sum_{A_j \in \mathcal{A}} \sum_{A \in \mathcal{A} \in A_j} R(A) = n \times n_p \times \log_2 K, \quad (7)$$

where  $n_p * \log_2 K$  denotes the constant cost of each image  $I \in \mathcal{I}$ . The bitrate  $R(\mathcal{Z}^m)$  and  $R(\mathcal{Z}^r)$  should be measured as the entropy of the corresponding latent representation symbols. In this paper, we use the entropy model in [15] to estimate  $R(\mathcal{Z}^m)$  and  $R(\mathcal{Z}^r)$ , respectively, which can be expressed as follows,

$$\begin{aligned} R(\mathcal{Z}^m) &= -\log_2 p(\mathcal{Z}^m | \mathcal{I}, \delta_m), \\ R(\mathcal{Z}^r) &= -\log_2 p(\mathcal{Z}^r | \delta_r), \end{aligned} \quad (8)$$

where  $\delta_m$  and  $\delta_r$  denote the parameters of the parametric and non-parametric entropy model, respectively, discriminated by whether they directly depend on the input images.  $p(\mathcal{Z}^m | \mathcal{I}, \delta_m)$  and  $p(\mathcal{Z}^r | \delta_r)$  denote the corresponding estimated probability of  $\mathcal{Z}^m$  and  $\mathcal{Z}^r$ .

### B. Gradient Approximation for Clustering Operation

Note that the Seb-based reference image generator employs a clustering algorithm, the module before the clustering algorithm (i.e., the Seb generator) cannot be updated in the back-propagation. Inspired by [19], in this paper, we directly copy the gradient of each of the reference patch generator's input  $S_i \in \mathcal{S}$  back to the corresponding outputs of the Seb generator  $\mathcal{F}_i = \{F_i | f_{\text{cluster}}(F_i) = S_i, F_i \in \mathcal{F}\}$ . However, the direct copy operation makes the clustering operation not constrained by E2E loss, resulting in an arbitrarily grown latent space. To tackle the problem, the  $L_2$  regulation is utilized to move each  $F_i$  towards  $S_i$ ,

$$L_{Reg} = Reg(\mathcal{F}, \text{sg}(\mathcal{S})) = \sum_{F_i \in \mathcal{F}} \sum_{S_i \in \mathcal{S}} \|F_i - \text{sg}(S_i)\|_2^2, \quad (9)$$

where  $\text{sg}(\cdot)$  represents the stop gradient operator that constraining  $\mathcal{S}$  to not be directly moved.

As a result, the total loss function can be expressed as

$$L = L_{RD} + \beta L_{Reg}, \quad (10)$$

where  $\beta$  controls the weight of the regulation loss. We use  $\beta = 1$  in our experiments.

## IV. EXPERIMENTAL RESULTS

In this section, we conduct extensive experiments to validate the effectiveness and generality of the proposed framework. The validation methodology is first provided, including the construction of training and testing datasets, the parameters setting, the choice of baselines, and the performance metrics for comparison. Then, the experimental results are further discussed to illustrate the performance of the proposed method.

### A. Methodology

As for training, a training set consisting of 50,000 images is sampled from the ImageNet training dataset [20], with each image being enhanced (randomly resized and cropped) into 16 augmentations with  $256 \times 256$  resolution as image subsets with a certain level of correlation during training. The Adam optimizer [21] is used by setting  $\beta_1$  and  $\beta_2$  as 0.9 and 0.999, respectively. The framework is trained at the learning rate of  $10^{-4}$  and  $10^{-5}$  each for one epoch, and continued at the learning rate of  $10^{-6}$  for three epochs until convergence.

To illustrate the advantages of the proposed framework, during validation, we construct a Mixed dataset by mixing samples from the Cityscapes' test set [22] and the UVG dataset [23]. These two datasets are composed of photos of urban street scenes and frames of video sequences, respectively, containing samples with distinct distributions. Specifically, the Mixed dataset consists of 500 images with 100 sampled from an UVG sequence, others sampled from Cityscapes' test set, and all images are cropped into  $1024 \times 1920$  resolution. The Mixed dataset denotes a representative use case of the proposed framework, where images with different levels of correlations and characteristics need to be transmitted.

We specify the Seb representation efficiency parameter  $K = \lfloor \frac{n_j * n_p}{25} \rfloor$  to make the number of Sebs with fixed proportion to the number of images in each subset, where  $\lfloor \cdot \rfloor$  denotes the floor operation. The patch height and width are set as  $h = w = 32$ . We compare the proposed framework with the following baselines,

- The wide-used engineered image compression codec JPEG2000 [24] as the representative baseline for traditional methods.
- The DL-based compression scheme [15] (denoted as "DL w/o corr."). The scheme utilizes an autoencoder structure for the compression and reconstruction of images, which each image be processed independently. The correlations among images are ignored.
- The DL-based compression scheme [16] (denoted as "DL w/ corr."). The scheme achieves compression based on image pairs. A prediction of the target image is made

based on the reference image, and then the autoencoders are used to compress the prediction and the residual parts, respectively. The correlations among images are strict.

Note that the same autoencoder structure is used in [15], [16] and the proposed framework. Therefore, the comparison with these two baselines avoids the impact caused by the difference in basic network structure to a certain extent, which can reflect the effectiveness of the proposed Seb-based method more accurately.

Without loss of generality, the experiments are taken over the AWGN channel, with each scheme combined with an ideal capacity-achieving channel code for transmission. We test the performance of each scheme under different signal-to-noise ratio (SNR) and channel bandwidth ratio (CBR) [25], [26] conditions, which reflect the channel condition and the overall coding rate, respectively. In specific, CBR is defined as the ratio between the number of channel input symbols (e.g. the number of transmitted symbols of  $\mathcal{S}$ ,  $\mathcal{A}$ ,  $\mathcal{Z}^m$ , and  $\mathcal{Z}^r$  for the proposed framework) and the number of source image symbols (e.g.  $3 \times 1024 \times 1920$  for images in  $\mathbb{R}^{3 \times 1024 \times 1920}$ ). In this case, we can derive the function of SNR and bit-per-pixel (BPP) of the above schemes under a given CBR, which by first the channel capacity  $C$  is obtained as

$$C = \frac{\text{BPP}}{3 \times \text{CBR}}, \quad (11)$$

and then SNR is obtained according to the Shannon channel capacity formula,

$$C = \log_2(1 + \text{SNR}). \quad (12)$$

Peak signal-to-noise ratio (PSNR) and multi-scale structural similarity (MS-SSIM) [27] are used to measure the quality of image reconstruction.

## B. Results and Analysis

Fig. 4 shows the PSNR and MS-SSIM results under different SNR conditions on the mixed dataset with CBR = 1/30. The proposed framework outperforms the baselines in general, achieving about 0.5 – 1.5 dB gain when measured in PSNR. Notably, the scheme [16] (DL w/ corr.) shows the worst performance due to the statistical difference between the training and testing data, which makes it unable to perform effective correlation information extraction. For the proposed framework, the performance gain is more significant under low SNR conditions, whereas the other DL-based schemes show more severe performance deterioration. This is because of the introduction of the Seb-based image representation mechanism, which supports an efficient recovery of the images' semantics while using limited communication resources. Besides, the framework also achieves similar or better performance w.r.t. MS-SSIM, as shown in Fig. 4(b).

Fig. 5 shows the average CBR consumption of the proposed framework corresponding to Sebs, usage information, and residual information under different  $\lambda$  ( $\lambda = \{256, 512, 1024, 2048\}$ ) and datasets with SNR = 5. It can be

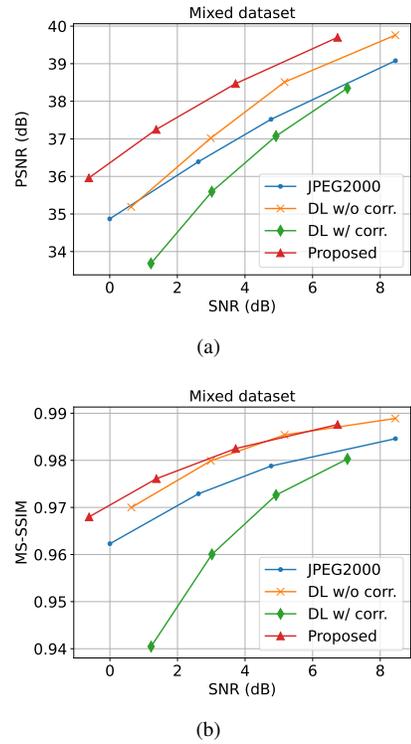


Fig. 4. SNR-distortion curves on the Mixed dataset with CBR = 1/30.

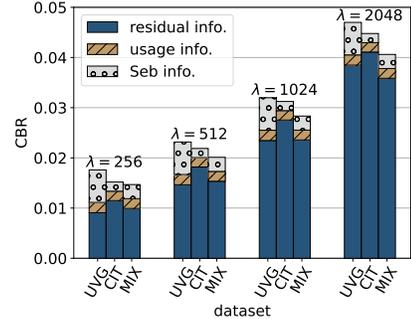


Fig. 5. The average CBR consumption of Sebs, usage information, and residual information for one image under different  $\lambda$  on the UVG, Cityscapes, and the Mixed datasets with SNR = 5.

observed that the proportion of Seb information varies significantly under different settings. The UVG dataset consumes the highest proportion of Seb information, followed by the Mixed and the Cityscapes datasets, respectively, under the same  $\lambda$  conditions. The result is in accordance with the correlation level of datasets that Seb carries most of the information for a strongly correlated dataset, leaving a small amount of residual information. For frameworks under different  $\lambda$ , a relatively fixed amount of Seb information, which depends on the characteristics (e.g. the complexity of image textures) of datasets, is consumed to support the recovery of images' semantics, and more residual information is consumed by frameworks with higher  $\lambda$  to satisfy the requirement of higher quality of image recovery. In addition, the CBR consumption of usage information is consistent with the result in (7), which

is nearly invariant under specific  $n_h$ ,  $n_w$ , and  $K$  settings.

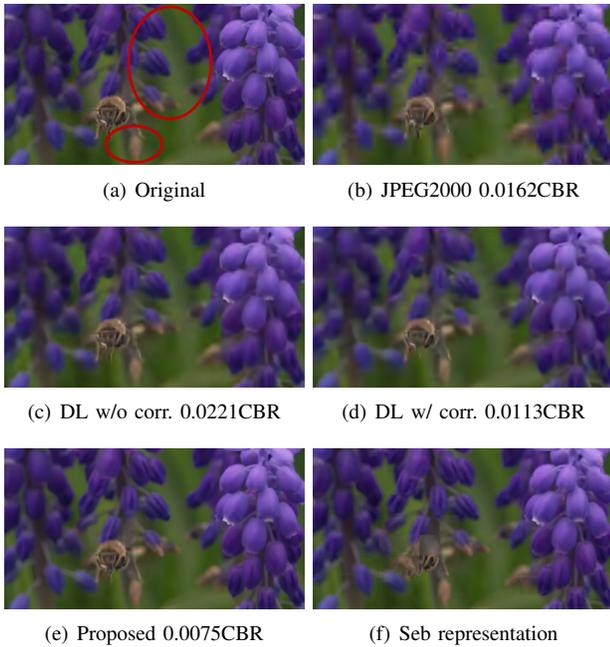


Fig. 6. Examples of reconstructions by different transmission schemes for low CBR values with SNR = 5.

Fig. 6 shows the reconstruction samples of different schemes. Compared to the relatively dynamic region (honey-bee), the static region (flower) is more precisely represented by Seb as shown in Fig. 6(f). The reason is that the static patterns are more densely distributed when mapping into the latent space, thus allowing for a better clustering result. Moreover, compared to the baselines, the proposed framework presents higher reconstruction quality (more precise color and details as shown in the oval region) with a significantly lower amount of transmitted data. Since the overall region of the image is effectively reconstructed by Sebs, only a small amount of residual is required to restore details, which is a typical use case for the proposed framework.

## V. CONCLUSION

In this paper, we propose a Seb-based image transmission framework, where common knowledge between the transmitter and the receiver is explicitly formed and shared in the form of Sebs. The framework includes a Seb-based reference image generator for Seb generation and Seb-based image representation, and a Seb-based image encoder/decoder to encode/decode the residual information for precise image reconstruction. A specialized loss function is introduced to solve the non-derivable problem. Experimental results show that the proposed framework outperforms baselines under all tested channel conditions. Future work would focus on further refining the generation and representation mechanism of Sebs.

## REFERENCES

[1] W. Jiang, B. Han, M. A. Habibi, and H. D. Schotten, "The road towards 6G: A comprehensive survey," *IEEE Open Journal of the Communications Society*, 334-366, 2021.

[2] K. Niu, J. Dai, and P. Zhang, "Semantic communication for 6G," *Mobile Communications*, vol. 45, no. 04, pp. 85-90, 2021.

[3] J. Liu, W. Zhang, and H. V. Poor, "A rate-distortion framework for characterizing semantic information," *IEEE International Symposium on Information Theory (ISIT)*, pp. 2894-2899, 2021.

[4] H. Xie, Z. Qin, G. Y. Li, and B. H. Juang, "Deep learning enabled semantic communication systems," *IEEE Transactions on Signal Processing* 69: 2663-2675, 2021.

[5] Q. Hu, G. Zhang, Z. Qin *et al.*, "Robust Semantic Communications Against Semantic Noise," *IEEE 96th Vehicular Technology Conference (VTC2022-Fall)*, pp. 1-6, 2022.

[6] Z. Qin, X. Tao, J. Lu, W. Tong, and G. Y. Li, "Semantic communications: Principles and challenges," *arXiv preprint arXiv:2201.01389*, 2021.

[7] G. Shi, D. Gao, X. Song *et al.*, "A new communication paradigm: from bit accuracy to semantic fidelity," *arXiv preprint arXiv:2101.12649*, 2021.

[8] W. Xu, Y. Zhang, F. Wang *et al.*, "Semantic Communication for the Internet of Vehicles: A Multiuser Cooperative Approach," *IEEE Vehicular Technology Magazine*, 2023.

[9] P. Jiang, C. K. Wen, S. Jin, and G. Y. Li, "Deep source-channel coding for sentence semantic transmission with HARQ," *IEEE Transactions on Communications*, 70.8: 5225-5240, 2022.

[10] H. Wei, W. Xu, F. Wang *et al.*, "SemAudio: Semantic-Aware Streaming Communications for Real-Time Audio Transmission," *IEEE Global Communications Conference*, 3965-3970, 2022.

[11] J. Hu, F. Wang, W. Xu, H. Gao, and P. Zhang, "Scalable Multi-Task Semantic Communication System with Feature Importance Ranking," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1-5, 2023.

[12] P. Zhang, W. Xu, H. Gao *et al.*, "Toward wisdom-evolutionary and primitive-concise 6G: A new paradigm of semantic communication networks," *Engineering*, 8: 60-73, 2022.

[13] X. Chen, and K. He, "Exploring simple siamese representation learning," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 15750-15758, 2021.

[14] J. Ballé, L. Valero, and E. P. Simoncelli, "Density Modeling of Images Using a Generalized Normalization Transformation," *4th Int. Conf. on Learning Representations*, 2016.

[15] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," *International Conference on Learning Representations*, 2018.

[16] G. Lu, W. Ouyang, D. Xu *et al.*, "Dvc: An end-to-end deep video compression framework," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.

[17] A. Ranjan and M. J. Black, "Optical flow estimation using a spatial pyramid network," *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4161-4170, 2017.

[18] F. Bellard, "BPG image format," URL: <https://bellard.org/bpg/>, 2015.

[19] A. Oord, and O. Vinyals, "Neural Discrete Representation Learning," *Advances in neural information processing systems*, 30, 2017.

[20] O. Russakovsky, J. Deng, H. Su *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, 115: 211-252, 2015.

[21] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[22] M. Cordts, M. Omran, S. Ramos *et al.*, "The cityscapes dataset for semantic urban scene understanding," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3213-3223, 2016.

[23] A. Mercat, M. Viitanen, and J. Vanne, "UVG dataset: 50/120fps 4K sequences for video codec analysis and development," in *Proc. ACM Multimedia Syst. Conf.*, Istanbul, Turkey, June 2020.

[24] JPEG2000, <https://www.openjpeg.org/>

[25] E. Boursoulatzé, 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.

[26] J. Wang, S. Wang, J. Dai *et al.*, "Perceptual learned source-channel coding for high-fidelity image semantic transmission," *IEEE Global Communications Conference*, 3959-3964, 2022.

[27] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers*, 2: 1398-1402, 2003.