

# Task-oriented Explainable Semantic Communications

Shuai Ma, Weining Qiao, Youlong Wu, Hang Li, Guangming Shi, *Fellow, IEEE*,  
Dahua Gao, Yuanming Shi, Shiyin Li, and Naofal Al-Dhahir, *Fellow, IEEE*

## Abstract

Semantic communications utilize the transceiver computing resources to alleviate scarce transmission resources, such as bandwidth and energy. Although the conventional deep learning (DL) based designs may achieve certain transmission efficiency, the uninterpretability issue of extracted features is the major challenge in the development of semantic communications. In this paper, we propose an explainable and robust semantic communication framework by incorporating the well-established bit-level communication system, which not only extracts and disentangles features into independent and semantically interpretable features, but also only selects task-relevant features for transmission, instead of all extracted features. Based on this framework, we derive the optimal input for rate-distortion-perception theory, and derive both lower and upper bounds on the semantic channel capacity. Furthermore, based on the  $\beta$ -variational autoencoder ( $\beta$ -VAE), we propose a practical explainable semantic communication system design, which simultaneously achieves semantic features selection and is robust against semantic channel noise. We further design a real-time wireless mobile semantic communication proof-of-concept prototype. Our simulations and experiments demonstrate that our proposed explainable semantic communications system can significantly improve transmission efficiency, and also verify the effectiveness of our proposed robust semantic transmission scheme.

## Index Terms

Explainable semantic communications, feature selection, semantic communications prototype

## I. INTRODUCTION

With the advent of augmented reality (AR), virtual reality (VR), holographic communications, autonomous vehicular networks, and industrial Internet of Things (IIoT), it is envisioned that

Shuai Ma is with Pengcheng Laboratory, Shenzhen, 518066, China (e-mail: mash01@pcl.ac.cn).

existing networks may soon reach a resource bottleneck due to stringent requirements [1], [2], such as ultra-high data rate, ultra-reliability, and low latency. To meet the above-mentioned requirements, investigations on the sixth generation communications (6G) are well underway and promise more powerful capacities than the fifth-generation communications (5G) [3]. From the first generation communications (1G) to 5G, the communication networks primarily focus on finding new resources and technologies to expand the channel capacity [4]. One approach is to seek the usage of large bandwidth, such as terahertz (THz) communications and visible light communication (VLC). Another approach is to explore the spatial domain, like ultra-massive MIMO and intelligent metasurfaces. However, given the hardware and physical limitations, the channel capacity may not keep increasing at the rate we desire to satisfy the aforementioned beyond-5G applications [5], [6].

In recent years, semantic communications, in which only task-relevant information is extracted and transmitted to the receiver, have received increasing attention by both academia and the industry [7]–[13]. Rather than increasing the channel capacity as in the conventional techniques, semantic communications exploit the computing power at the transceivers to alleviate the cost of transmission resources. The classic Shannon information theory focuses on “How accurately can the symbols be transmitted?”, which ignores the meaning of the transmitted messages. Instead, semantic communications [14] consider “How precisely do the transmitted symbols convey the desired meaning?” Thus, it is possible to improve the system efficiency at the semantic level, not only at the pure bit level.

The classical separation theorem [15] states that, as the data size goes to infinity, separating source coding and channel coding can achieve the optimal performance over a memoryless communication channel. However, for finite number of bits transmission, the performance of such separated structure will degrade. This issue also arises in semantic communications. Various deep learning (DL) based joint source-channel coding (JSCC) schemes have been investigated for text [12], [16], [17], image [18]–[23], speech [24], [25], and multimodal data [26] transmission. Specifically, for text semantic transmission, the JSCC schemes have been designed by exploiting architectures like the recurrent neural network (RNN) [16], Transformer [17], [27], autoencoder (AE) [28], adaptive Universal Transformer [29], and deep neural network (DNN) [30]. For image semantic transmission, a masked auto-encoder (MAE) architecture with Transformer was designed in [18] to combat adversarial samples noise. Convolutional neural networks (CNNs) based JSCC schemes were designed for the time-invariant and fading wireless channels in [19].

Neural error correcting and source trimming (NECST) codes were studied in [22]. For finite bit transmission, an attention DL based JSCC method was designed in [23]. By exploring the channel output feedback, an AE-based JSCC scheme was developed in [20] to improve the quality of image transmission. By combining an AE with orthogonal frequency division multiplexing (OFDM), a JSCC wireless image transmission scheme was presented in [21] over multipath fading channels. By leveraging reinforcement learning (RL), a joint semantics-noise coding (JSNC) mechanism was designed in [31]. A DNN based JSCC scheme was designed in [32] for adaptive rate control in wireless image transmission. Based on AE, a SNR-adaptive deep JSCC scheme is proposed in [33] for multi-user wireless image transmission. To tackle the variational information bottleneck, the authors in [34] investigated task-oriented communication for edge inference, where a low-end edge device extracts the feature vector of a local data sample and transmits to a powerful edge server for processing. Besides, for the speech semantic transmission, AE based wave-to-vector architecture and squeeze-and-excitation (SE) attention network have been developed in [24] and [25], respectively. For visual question answering, the memory-attention-composition neural network was designed in [26] for multi-modal data semantic communications.

However, most of the existing works on semantic communications [16]–[26] are based on DL techniques, in which the DL model is basically a black box. Thus, the extracted semantic feature vectors in these works are unexplainable (hidden) representations, and the uninterpretability of the extracted features restricts further processing and exploitation of semantic features. For example, due to the uninterpretability, the unintended features will also be transmitted to the receiver, which wastes transmission resources and reduces the efficiency of semantic communications.

Moreover, most of the existing semantic communication investigations [16]–[26] completely redesign the source and channel module over the conventional system, which are impractical and not compatible with the existing communication networks. Because there is a large number of practical standards and hardware for 5G physical layer, it will lead to a huge waste of resources and costs by replacing physical layer techniques with DL-based semantic JSCC techniques. Therefore, how to design efficient and 5G-compatible practical semantic communications is a critical issue.

To address the above the two key challenges of the semantic communications, we propose an explainable and robust semantic communication framework in this paper, which is compatible with existing communication systems. We show that the proposed framework can achieve a

higher transmission efficiency than the existing inexplicable semantic communication systems. The main contributions of this paper are summarized as follows:

- We propose an explainable and easy-to-implement semantic communication framework based on the bit-level communication systems, which includes a novel semantic encoder, as well as the corresponding decoder, feature selection and semantic channel. The innovation of the proposed framework is threefold: i) The semantic encoder/decoder aims to, not only extract the independent and explainable semantic information as *semantic source coding*, but also alleviate the ambiguity of the semantic information influenced by the quantization and channel noise as *semantic channel coding*; ii) The feature selection module follows the semantic encoder, to choose only the task-relevant features for transmission, which can further reduce the transmission load; iii) The framework has an explicit definition of semantic channels, which incorporates the key modules of the bit-level communication systems. Specifically, the semantic channel takes both quantization error (or noise) and physical channel noise into account since those noise sources may lead to semantic information ambiguity, and the semantic channel capsulizes the conventional bit-level communication systems, which implies that the proposed framework can be more easily implemented compared to the JSCC schemes.
- Then, we propose two information-theoretic metrics for our semantic communication framework. In terms of the information compression of the semantic encoder, we derive the optimal distribution of the reconstruction signal of the rate-distortion-perception function for semantic information extraction. Moreover, to quantify the semantic information transmission, we derive both upper and lower bounds for the semantic channel capacity, which are shown to be tight when the quantization noise tends to zero.
- Based on our framework, we further propose a feasible design of the explainable semantic communication system. Specifically, this design includes a robust  $\beta$ -VAE lightweight unsupervised learning network, where a weighted parameter  $\beta$  is added to the Kullback-Leibler (KL) divergence term of the variational autoencoder (VAE) network loss function, in order to make the latent representations effectively disentangled. Moreover, to enhance transmission robustness, the semantic channel noise is added to the extracted features during semantic networks training.
- Finally, we implement the above semantic communication design, and propose a wireless

mobile semantic communication proof-of-concept prototype. Applying the portable Raspberry Pi 4 Model B and Wi-Fi, the developed prototype can run the proposed robust  $\beta$ -VAE semantic system in real time. Our experiments demonstrate that our proposed semantic communication system can achieve better performance than existing benchmarks.

The rest of this paper is organized as follows. The explainable semantic communications framework is presented in Section II. Section III provides the information-theoretic metrics of semantic communications. In Section IV, we propose a  $\beta$ -VAE based robust and explainable semantic communications system. In Section V, we present the semantic communication system prototype design and implementation. In Section VI, we evaluate the proposed explainable semantic communication system. Finally, we conclude the paper in Section VII. Table I and II presents the means of the key notations and key acronyms in this paper, respectively.

TABLE I: Key Notations and Meanings

| Variables  | Meanings  |
|--|---|
| $S = \{s_k\}_{k=1}^K$                                  | Semantic information with $K$ features                  |
| $s_k$  | The $k$ th semantic feature                             |
| $X$  | Source data   |
| $Z = \{z_l\}_{l=1}^L$                                  | The extracted semantic feature vector with $L$ features |
| $z_l$  | The $l$ -th extracted semantic feature                  |
| $\mathcal{L}$  | Semantic feature index set                              |
| $\mathcal{L}_{\text{sel}}$                             | Selected semantic feature index set                     |
| $X_s = \{z_l\}_{l \in \mathcal{L}_{\text{sel}}}$       | Selected semantic features                              |
| $Y_s = \{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}}$ | Estimated semantic features                             |
| $\hat{Z}$  | Reconstructed feature set                               |
| $\hat{X}$  | Decoded data  |

TABLE II: Key Acronyms and Meanings

| Acronyms | Meanings                            |
|----------|-------------------------------------|
| JSCC     | Joint source-channel coding         |
| VAE      | Variational autoencoder             |
| KL       | Kullback- Leibler                   |
| ANGC     | Additive non-Gaussian noise channel |
| ELBO     | Evidence lower bound                |
| GPU      | Graphics processing unit            |
| PSNR     | Peak signal-to-noise ratio          |

## II. EXPLAINABLE SEMANTIC COMMUNICATION FRAMEWORK

Most existing studies replace the traditional source coding and channel coding modules by deep learning-based studies source-channel coding, which greatly changes the structure of the existing communication systems. In this paper, we propose a semantic communication framework incorporating the key modules of the conventional communication system (e.g., 5G).

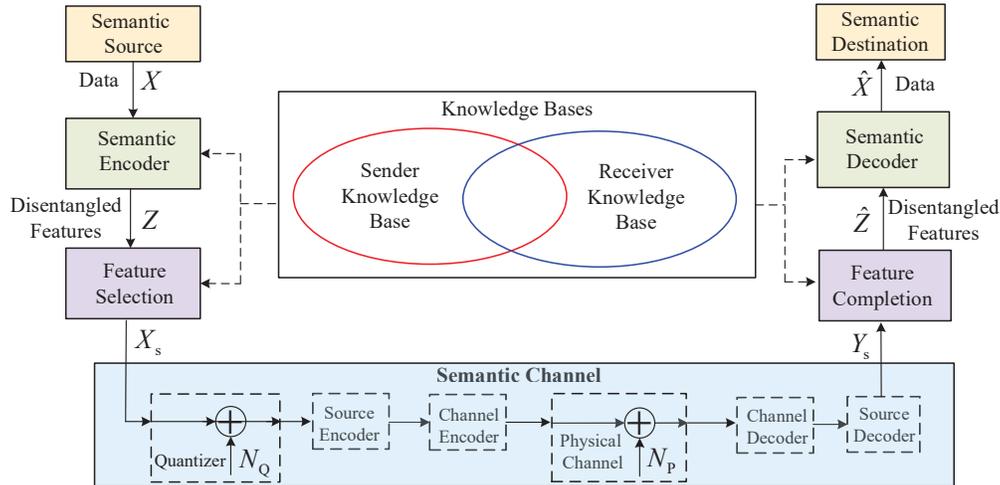


Fig. 1: Explainable and robust semantic communications framework

As shown in Fig. 1, the proposed explainable semantic communication framework includes a semantic source, sender knowledge base, semantic encoder, semantic channel, receiver knowledge base, semantic decoder, and semantic destination. Note that the proposed framework introduces a semantic-level transmission on the top of bit-level transmission. Clearly, such a framework does not require the extra redesign over the existing physical standards, protocols and products, which makes the application of semantic communications more practical. Next, we will describe each module in detail.

### A. Knowledge Bases

The knowledge base contains all the necessary information that can facilitate the communication at the semantic level. Specifically, the knowledge base includes background knowledge and training dataset. The background knowledge is used to facilitate the semantic feature extraction and selection in the semantic transmitter. The training dataset is used for training the parameters of the semantic encoder and decoder. The sender may choose different semantic knowledge bases according to different tasks, scenarios and recipients. For example, when the communication is

triggered between people in different countries, it may be necessary to sample multiple language databases. In general, the sender and the receiver share some common knowledge, which may act as a special kind of side information to improve coding efficiency.

### B. Semantic Sources

The semantic source produces original data, such as pictures, videos, voices, and texts. The generated data contains certain semantic information to be shared with the semantic destination. Assume that the semantic information includes  $K$  features  $S = \{s_k\}_{k=1}^K \sim p_{\text{sou}}(s)$  (data generative factors), where  $s_k$  denotes the  $k$ th semantic feature, and the joint probability distribution is  $p_{\text{sou}}(s)$ . Further, assume that the  $K$  features  $\{s_k\}_{k=1}^K$  are independent, i.e.,  $p_{\text{sou}}(s) = \prod_{k=1}^K p_{\text{sou}}(s_k)$ , where  $p_{\text{sou}}(s_k)$  denotes the probability distribution of  $s_k$ . Thus, the entropy of the semantic source is given as

$$H(S) = - \sum_{k=1}^K p_{\text{sou}}(s_k) \log_2 p_{\text{sou}}(s_k). \quad (1)$$

The semantic source can generate the data  $X \sim p_{\text{data}}(x)$ , which can be images, text, sound or video. Generally, the generated data need to include both the intended features and some redundant features to make the whole semantic data complete. Thus, the data generation is defined as  $p_{\text{s2d}}(x | \{s_k\}_{k=1}^K)$ , and the probability distribution function (PDF) of data  $X$  is given as

$$p_{\text{data}}(x) = \sum_{s_1, \dots, s_K} p_{\text{s2d}}(x | \{s_k\}_{k=1}^K) \prod_{k=1}^K p_{\text{sou}}(s_k). \quad (2)$$

The entropy of the semantic data  $x$  is given as

$$H(X) = - \sum_x p_{\text{data}}(x) \log_2 p_{\text{data}}(x). \quad (3)$$

Based on (1) and (2),  $H(X)$  can be further expressed as

$$H(X) = H(S) + H(X|S) - H(S|X). \quad (4)$$

### C. Semantic Encoder

Based on the knowledge base, the generated message  $X$  will be processed by the semantic encoder, which is a joint semantic source and channel encoder. More specifically, it extracts

semantic information or the semantic features of the message  $X$ , and outputs the disentangled and explainable features  $Z$ , which can be viewed as a semantic source encoder. On the other hand, in order to reduce the ambiguity incurred by the quantization error and channel noise, the semantic encoder needs to improve the robustness against the semantic channel noise, which can be viewed as a semantic channel encoder.

The semantic encoder extracts a low-dimensional semantic features vector  $Z \sim p_{\text{fea}}(z)$  from the data  $X$ . Let  $p_{\text{d2f}}(z|x)$  denote the conditional PDF of the feature  $z$  given data  $x$ . Thus, the PDF of the extracted feature (sub-vectors)  $p_{\text{fea}}(z)$  is given as

$$p_{\text{fea}}(z) = \sum_x p_{\text{d2f}}(z|x) p_{\text{data}}(x). \quad (5)$$

The encoder is required to regulate the extracted features into  $L$  independent features  $Z = \{z_l\}_{l=1}^L$ , which satisfy

$$p_{\text{fea}}(z) = \prod_{l=1}^L p_{\text{fea}}(z_l), \quad (6)$$

where  $p(z_l)$  denote the PDF for feature  $z_l$ . In summary, the extracted feature  $\mathbf{z}$  is required to have  $L$  disentangled interpretable semantic features  $\{z_l\}_{l=1}^L$ , whose corresponding neural network output is explainable and understandable by the human. For convenience, we let  $\mathcal{L} \triangleq \{1, \dots, L\}$  denote the index set of the disentangled semantic features. Note that such a requirement can be met if the semantic encoder is designed in a sophisticated manner. In Section IV, we will introduce a feasible system design that has such capability.

#### D. Feature Selection

It should be noted that the obtained features  $Z$  could contain more information than what the receiver is interested in. Thus, after extracting the disentangled features  $\{z_l\}_{l=1}^L$ , only the subset of features  $\{z_l\}_{l=1}^L$  that are of interest to the receiver should be transmitted, and the rest of the features can be viewed as the ‘‘redundancy’’. We will present more discussions of this issue via experiments in Section VI.

Given the task requirement, let  $\mathcal{L}_{\text{sel}} \subseteq \mathcal{L}$  denote the selected feature index set, then the selected set of features is given as

$$X_s = \{z_l\}_{l \in \mathcal{L}_{\text{sel}}}. \quad (7)$$

Thus, feature selection will reduce the amount of data sent, and the corresponding reduction is  $\{z_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$ . Then,  $X_s$  will be sent to the semantic channel.

### E. Semantic Channel

After the feature selection module, the task-oriented features are selected and ready to send. Since the quantization error and channel noise both could incur semantic information ambiguity, we define the semantic channel with channel law  $p(y_s|x_s)$  as a virtual channel including the signal quantizer and the bit-level communication system, as shown in Fig. 1. Here,  $Y_s = \{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}} \sim p_r(y_s)$  represents the set of estimated features after the transmissions over the semantic channel, and  $\hat{z}_l$  denotes the estimated feature of  $z_l$ .

Generally, the semantic noise could include various factors including source errors, feature extraction errors, knowledge base ambiguities, adversarial injections, quantization noise, physical channel noise, etc. In our framework, the semantic channel noise  $N_s$  is the distortion between the selected semantic feature  $X_s = \{z_l\}_{l \in \mathcal{L}_{\text{sel}}}$  and the estimated semantic feature  $Y_s = \{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}}$ , which mainly depends on the quantization noise and physical channel noise.

1) *Quantization Noise*: The quantization noise is caused by the traditional communication operation modules, such as the source encoder (or decoder) or channel encoder (or decoder), which may also lead to semantic ambiguity. In order to reduce the number of transmitted bits, the semantic feature  $\mathbf{x}_s$  will be converted to a compressible binary stream using few bits. To represent  $\mathbf{x}_s$  with a finite number of bits, we need to map it to a discrete space. Specifically, a finite quantizer maps the semantic feature  $x_s$  to  $x_b$ , whose values are then quantized to  $M$  levels  $C = \{c_1, \dots, c_M\}$ , i.e.,

$$x_b = \text{Quan}(x_s), \quad (8)$$

where  $\text{Quan}(\cdot)$  is a quantization operator. Since the number of dimensions  $\dim(x_b)$  and the number of levels  $L$  are finite, the entropy of quantized semantic data is given as

$$H(x_b) \leq \dim(x_b) \log_2 M. \quad (9)$$

In this paper, we consider the uniform distributed quantization noise  $n_Q$ , i.e.,

$$p_{N_Q}(x) = \frac{1}{b-a}, a \leq x \leq b, \quad (10)$$

where  $a$  and  $b$  are the lower and upper bounds of quantization noise  $N_Q$ .

2) *Physical Channel Noise*: The physical channel noise exists ubiquitously in physical communications and is caused by physical channel impairments, such as additive white Gaussian noise (AWGN), interference, etc. It is noted that the errors caused by channel propagation usually occur before channel decoding and can be corrected by channel decoding. Assume that the physical channel noise  $N_P$  follows a Gaussian distribution with zero-mean and variance  $\sigma_P^2$ , i.e.,

$$p_{N_P}(x) = \frac{1}{\sigma_P \sqrt{2\pi}} \exp\left(-\frac{x^2}{2\sigma_P^2}\right). \quad (11)$$

### F. Feature Completion

After obtaining the estimated features  $Y_s = \{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}}$  through the semantic channel transmission, the destination will use the estimated features and side information in the knowledge base, to compute the target function of the task. Although the unintended features subset  $\{z_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$  are not transmitted, the receiver may generate the corresponding unintended features  $\{\hat{z}_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$  by exploiting the knowledge base. Then, by combining intended features  $\{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}}$  and unintended features  $\{\hat{z}_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$ , we may obtain the completed semantic features  $\hat{Z} = \{\hat{z}_l\}_{l \in \mathcal{L}}$  with distribution  $p_{\text{rfea}}(\hat{z})$ . For example, considering a semantic communication system for staff clothing image transmission, the intended semantic features of the receiver are clothing features, and the receiver is not interested in the staff's gender, skin color, and hairstyle. Therefore, the receiver can generate unintended semantic features based on the shared knowledge base, such as the staff's gender, skin color and hairstyle. Note that, the generated unintended semantic features at the receiver may be different from the corresponding features of the image at the transmitter. Then, the receiver combines the received clothing features with its own generated unintended features.

### G. Semantic Decoder

The semantic decoder aims to recover the data from the disentangled features  $\hat{Z}$  that are semantic explainable, which is the inverse function of the semantic encoding. Again, this inverse function needs the help of the knowledge base for model training such that the decoder can “understand” the features  $\hat{Z}$ .

Similar to the encoding process, we use conditional PDF  $p_{f2d}(\hat{x}|\hat{z})$  to describe the semantic decoding process. The PDF of the decoded data is  $p_{\text{rdata}}(\hat{x})$ , and the decoded data is  $\hat{x}$ . The data reconstruction for a given feature vector is given as

$$p_{\text{rdata}}(\hat{x}) = \sum_{\hat{z}} p_{f2d}(\hat{x}|\hat{z}) p_{\text{rfea}}(\hat{z}). \quad (12)$$

#### H. Semantic Destination

Finally, the receiver recovers the semantic information based on the decoded data  $\hat{X}$ , and the corresponding process can be described by  $p_{\text{d2s}}(\hat{s}|\hat{x})$ , where the final semantic information is denoted by  $\hat{s}$ . At last, the probability of such semantic information can be written as

$$p_{\text{des}}(\hat{s}) = p_{\text{d2s}}(\hat{s}|\hat{x}) p_{\text{rdata}}(\hat{x}). \quad (13)$$

So far, we have presented the complete semantic communication framework. The key modules are the semantic encoder and the feature selection. Their functions can be realized by the careful model design. We will present a detailed system design in Section IV, which is a feasible realization of this framework.

### III. INFORMATION-THEORETIC METRICS OF SEMANTIC COMMUNICATIONS

In this section, we propose two metrics for the framework illustrated by Fig. 1. Here, we focus on two procedures: the encoding and the transmission.

#### A. Rate-Distortion-Perception Function

The semantic encoding may include many different tasks, and these tasks may have relevant or different criteria. For example, there is data distortion for the traditional data reconstruction task, and distribution distortion for generative learning tasks.

Let  $p(x)$  be the distribution of the input source,  $r(\hat{x})$  be the distribution of the reconstruction signal, and  $q(\hat{x}|x)$  be a conditional distribution on  $\mathcal{X} \times \mathcal{X}$ . The information rate-distortion-perception function  $R(D, P)$  [35] for a source  $X \sim p(x)$  is defined as

$$R(D, P) = \min_{q(\hat{x}|x)} I(X; \hat{X}) \quad (14a)$$

$$\text{s.t. } \mathbb{E}[\Delta(x, \hat{x})] \leq D, \quad (14b)$$

$$d(p(x), r(\hat{x})) \leq P, \quad (14c)$$

$$\sum_{\hat{x}} q(\hat{x}|x) = 1, \forall x \in \mathcal{X}. \quad (14d)$$

where the distortion function  $\Delta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}^+$  satisfying  $\Delta(x, \hat{x}) = 0$  if  $x = \hat{x}$ , and perception function  $d(p(x), r(\hat{x}))$  is a non-negative divergence between probability distributions  $p(x)$  and  $r(\hat{x})$  satisfying  $d(p, q) = 0$  if  $p(x) = r(x)$ .

So far, for a general source, the optimal distribution of the reconstruction signal  $r(\hat{x})$  of problem (14) has not been derived yet. For a binary source, the three-way tradeoff between rate, distortion, and perception was investigated in [35] with Hamming distance distortion and total-variation distance perception. While for a Gaussian source, the achievable distortion-perception region was established in [36] under squared error distortion and squared Wasserstein-2 distance.

Hence, we investigate how to find the optimal of  $R(D, P)$  for a general source under the mean square distortion (i.e.,  $\Delta(x, \hat{x}) = |x - \hat{x}|^2$ ) and KL divergence perception (i.e.,  $d(p(x), r(\hat{x})) = d_{\text{KL}}(p(x), r(\hat{x})) \triangleq \sum_x p(x) \log \frac{p(x)}{r(\hat{x})}$ ). We first introduce the following lemma.

**Lemma 1.** *Consider the mean square distortion (i.e.,  $\Delta(x, \hat{x}) = |x - \hat{x}|^2$ ) and KL divergence perception (i.e.,  $d(p(x), r(\hat{x})) = \sum_x p(x) \log \frac{p(x)}{r(\hat{x})}$ ). The corresponding optimal distribution  $q^*(\hat{x}|x)$  to problem (14) for a given output distribution  $r(\hat{x}) > 0$  is*

$$q^*(\hat{x}|x) = \frac{r(\hat{x})}{\tilde{\gamma}(x)} \exp\left(\mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2\right), \quad (15)$$

where  $\tilde{\gamma}(x) = \sum_{\hat{x}} r(\hat{x}) \exp\left(\mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2\right)$ . The corresponding optimal distribution  $r^*(\hat{x})$  to (14) for a given conditional distribution  $q(\hat{x}|x) > 0$  is

$$r^*(\hat{x}) = \sum_x p(x) q(\hat{x}|x). \quad (16)$$

*Proof:* Please find the proof in Appendix A. □

Using Lemma 1, we can apply a process of alternating minimization, called the Blahut–Arimoto algorithm [37]. Specifically, in the initialization setup, choose some positive values  $\alpha, \mu$  and the initial output distribution  $r^{(0)}(\hat{x})$ . In each iteration  $k$ , compute the optimal  $q^{(k)}(\hat{x}|x)$  according to (15) for given  $r^{(k-1)}(x)$ , and then compute the optimal  $r^{(k)}(x)$  according to (16).

## B. Lower and Upper Bounds on Semantic Channel Capacity

The channel capacity quantifies the maximum rate of information transmission for the considered system. According to the framework in Fig. 1, we define the semantic channel capacity

as the maximum semantic information that can be transferred through the semantic channel  $p(y_s|x_s)$ . Following the standard achievability and converse proof techniques, we obtain the semantic channel capacity in our framework as:

$$C_s = \max_{p(x_s)} I(X_s; Y_s). \quad (17)$$

In the conventional bit-level wireless communication system, the channel capacity is usually represented by the Shannon capacity formula with additive Gaussian distributed noise. In our framework, the semantic channel noise  $N_s$  mainly depends on the quantization noise and physical channel noise, and follows non-Gaussian distribution in general. Thus, in our framework, the semantic channel is an additive non-Gaussian noise channel (ANGC), and we assume that the estimated semantic features  $Y_s$  can be represented as

$$Y_s = X_s + N_s. \quad (18)$$

Although the specific distribution of  $N_s$  is unknown, the variance of the semantic noise  $n_s$  can be obtained by measurement. In this paper, we assume that the covariance of the semantic noise  $n_s$  is  $\sigma_s^2$ .

Due to the non-Gaussian distributed semantic noise  $n_s$ , the classic Shannon capacity formula (based on Gaussian distributed noise) cannot be directly applied to the semantic channel. To derive the semantic channel capacity, we first define equivalent Gaussian distributed semantic channel noise  $\bar{N}_s \sim \mathcal{N}(0, \sigma_s^2)$  with the same variance as  $N_s$ . Then, based on the equivalent semantic channel noise  $\bar{N}_s$ , the received signal of semantic channel  $\bar{Y}_s$  is given as

$$\bar{Y}_s = X_s + \bar{N}_s. \quad (19)$$

Therefore, the channel capacity of the equivalent semantic channel is given as

$$C_{s,\text{eq}} = \frac{1}{2} \log \left( 1 + \frac{P_{x_s}}{\sigma_s^2} \right), \quad (20)$$

where  $P_{x_s}$  denote the power of transmitted semantic data  $X_s$ .

**Proposition 1** (Lower and upper bounds on the semantic channel capacity). *With the non-Gaussian distributed channel noise, the semantic channel capacity  $C_s$  is bounded by [38]*

$$C_{s,\text{eq}} \leq C_s \leq C_{s,\text{eq}} + d_{\text{KL}}(p_{n_s}(x), p_{\bar{n}_s}(x)), \quad (21)$$

where  $d_{\text{KL}}(p_{n_s}(x), p_{\bar{n}_s}(x)) = \int_{-\infty}^{\infty} p_{n_s}(x) \log \frac{p_{n_s}(x)}{p_{\bar{n}_s}(x)} dx$ .

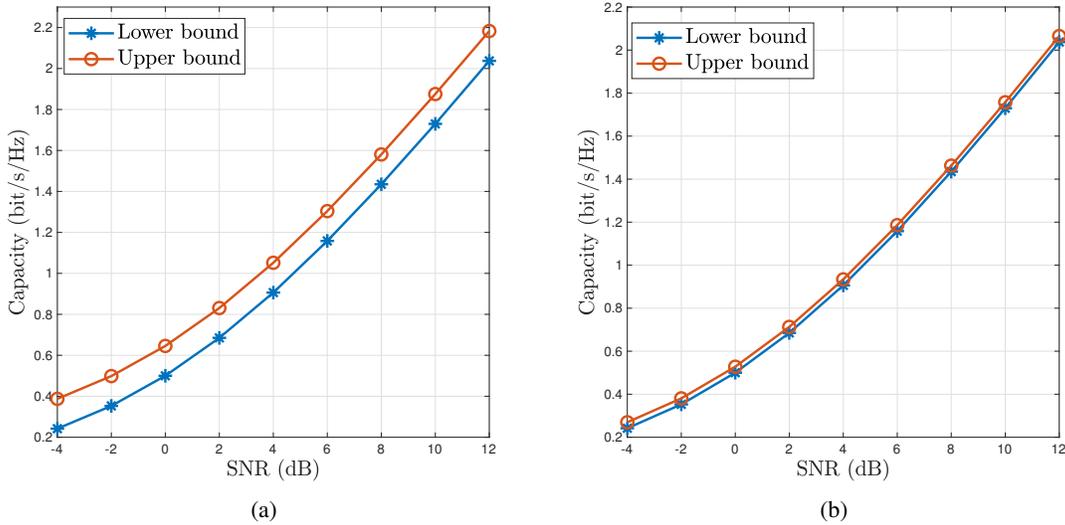


Fig. 2: (a) Lower and upper bounds of semantic channel capacity versus SNR with  $a = -1$ ,  $b = 1$  and  $\sigma_P^2 = 0.01$ ; (b) Lower and upper bounds of semantic channel capacity versus SNR with  $a = -0.3$ ,  $b = 0.3$  and  $\sigma_P^2 = 0.01$ .

At last, we illustrate our theoretical results on the semantic channel capacity via numerical simulation. Fig. 2 (a) and (b) show the lower bound and the upper bound in (22) on semantic channel capacity versus SNR with semantic noise parameters  $a = -1$ ,  $b = 1$  and  $\sigma_P^2 = 0.01$ , and semantic noise parameters  $a = -0.3$ ,  $b = 0.3$  and  $\sigma_P^2 = 0.01$ , respectively. Fig. 2 (b) shows that the gap between the upper bound and lower bound is less than that in Fig. 2 (a). The reason is that when the the KL divergence between semantic noise  $\bar{N}_s$  and the equivalent semantic channel noise  $\bar{n}_s$  tends to 0, i.e.,  $d_{\text{KL}}(p_{n_s}(x), p_{\bar{n}_s}(x)) \rightarrow 0$ , the gap between the lower bound and the upper bound in (22) tends to 0.

#### IV. $\beta$ -VAE BASED ROBUST AND EXPLAINABLE SEMANTIC COMMUNICATION SYSTEM

In this section, we present a feasible and efficient system design based on the proposed framework given in Fig. 1. Here, we propose a robust  $\beta$ -VAE based semantic communications system, as shown in Fig. 3, which disentangles the hidden representation vector into multiple independent and semantically interpretable of features.

##### A. Robust $\beta$ -VAE based Semantic Encoder/Decoder

By exploiting a generative VAE model [39], we first optimize the semantic encoder  $q_\phi(z|x)$  with parameter set  $\phi$ , and the semantic decoder  $p_\theta(\hat{x}|\hat{z})$  for the receiver with parameter set  $\theta$ .

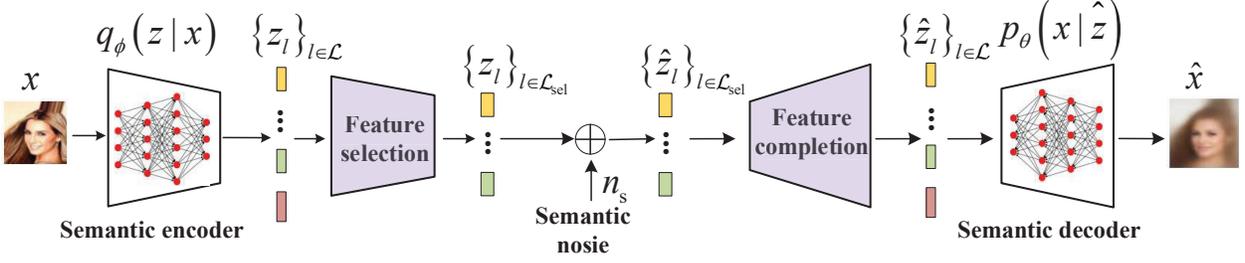


Fig. 3: Proposed  $\beta$ -VAE based explainable semantic communication system

Mathematically, we aim to jointly optimize parameters  $\phi$  and  $\theta$  to maximize the log-likelihood of data  $X$  as follows

$$\max_{\phi, \theta} \log p_{\theta}(x). \quad (22)$$

To efficiently handle optimization problem (22), we optimize the lower bound of the objective function  $\log p_{\theta}(x)$  [40]. Specifically,  $\log p_{\theta}(x)$  is lower bounded by

$$\log p_{\theta}(x) = \int_z q_{\phi}(z|x) \log p_{\theta}(x) dz \quad (23a)$$

$$= \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(z, x)}{p_{\theta}(z|x)} dz \quad (23b)$$

$$= \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(z, x)}{q_{\phi}(z|x)} dz + \int_z q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} dz \quad (23c)$$

$$= \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(z, x)}{q_{\phi}(z|x)} dz + d_{\text{KL}}(q_{\phi}(z|x) || p_{\theta}(z|x)) \quad (23d)$$

$$\geq \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(z, x)}{q_{\phi}(z|x)} dz, \quad (23e)$$

$$= \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(x|z) p_{\theta}(z)}{q_{\phi}(z|x)} dz \quad (23f)$$

$$= \int_z q_{\phi}(z|x) \log p_{\theta}(x|z) dz + \int_z q_{\phi}(z|x) \log \frac{p_{\theta}(z)}{q_{\phi}(z|x)} dz \quad (23g)$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - d_{\text{KL}}(q_{\phi}(z|x) || p_{\theta}(z)) \quad (23h)$$

where equation (23a) holds for the arbitrary distribution  $q_{\phi}(z|x)$ , and inequality (23e) holds due to  $d_{\text{KL}}(q_{\phi}(z|x) || p_{\theta}(z|x)) \geq 0$ .

Unfortunately, maximizing the lower bound in (23h) directly cannot achieve interpretable and robust semantic communication systems design. To address this challenge, we multiply  $d_{\text{KL}}(q_\phi(z|x)||p_\theta(z))$  by a weighting parameter  $\beta$  to obtain a disentangling and explainable semantic representation  $z$  [39], for  $\beta > 1$ . Furthermore, to combat semantic noise and achieve robust semantic communication systems design, we replace  $p_\theta(z|x)$  with  $p_\theta(x|\hat{z})$ , where  $\hat{z} = gz + n_s$ ,  $g$  denotes fading channel gain, and  $n_s$  denotes the semantic noise. Specifically, the log-likelihood maximization problem (22) is reformulated as follows

$$\max_{\phi, \theta} \mathbb{E}_{q_\phi(x|z)} [\log p_\theta(x|\hat{z})] - \beta d_{\text{KL}}(q_\phi(z|x)||p_\theta(z)), \quad (24)$$

where the prior distribution  $p_\theta(z)$  is assumed to follow a standard Gaussian distribution, i.e.,  $p_\theta(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Note that, in (24), the first term  $\mathbb{E}_{q_\phi(x|z)} [\log p_\theta(x|z + n_s)]$  is the expected likelihood with the cross entropy form, which can be regarded as reconstruction loss, while the second term regularizes  $q_\phi(z|x)$  to be close to prior  $p_\theta(z)$ , which can be regarded as regularization loss. To further enhance the robustness of the variational inference, we exploit  $\eta$ -cross entropy  $c_\eta(p(x)||p_\theta(x|\hat{z}))$  [41], [42] as the reconstruction loss, instead of the cross entropy  $\mathbb{E}_{q_\phi(x|z)} [\log p_\theta(x|z + n_s)]$ , where

$$c_\eta(p(x)||p_\theta(x|\hat{z})) = -\frac{\eta + 1}{\eta} \int p(x)^\eta dx + \int p_\theta(x|\hat{z})^{1+\eta} dx. \quad (25)$$

Specifically, the objective function of the proposed robust semantic communication system is given as

$$\max_{\phi, \theta} c_\eta(p(x)||p_\theta(x|\hat{z})) - \beta d_{\text{KL}}(q_\phi(z|x)||p_\theta(z)). \quad (26)$$

Thus, the robust  $\beta$ -VAE training objective (26) encourages the latent distribution  $q_\phi(z|x)$  to efficiently represent semantic information about the data  $x$  by jointly maximizing the  $\eta$ -cross entropy  $c_\eta(p(x)||p_\theta(x|\hat{z}))$  and minimizing the  $\beta$ -weighted KL term  $d_{\text{KL}}(q_\phi(z|x)||p_\theta(z))$  via unsupervised learning.

More specifically, we jointly optimize the semantic encoder parameter  $\phi$  and semantic decoder parameter  $\theta$  to maximize the objective function (26). The first term of (26) is the probability of reconstructing the input data  $x$ , which corresponds to reconstruction loss. The second term is minimizing the KL divergence, which is the distance between the approximated posterior  $q_\phi(z|x)$  and the fixed Gaussian distribution  $p_\theta(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ . By adopting the well chosen

values of the parameter  $\beta$  (usually  $\beta > 1$ ), the posterior  $q_\phi(z|x)$  is encouraged to match the Gaussian distribution  $p_\theta(z) = \mathcal{N}(0, \mathbf{I})$ , which disentangles the hidden representation into multiple independent and semantically meaningful features  $\{z_l\}_{l \in \mathcal{L}}$ . The parameter  $\beta$  balances reconstruction accuracy and learned disentanglement quality. In general, a higher value of  $\beta$  will produce a more disentangled representation, but may lead to lower reconstruction accuracy [39].

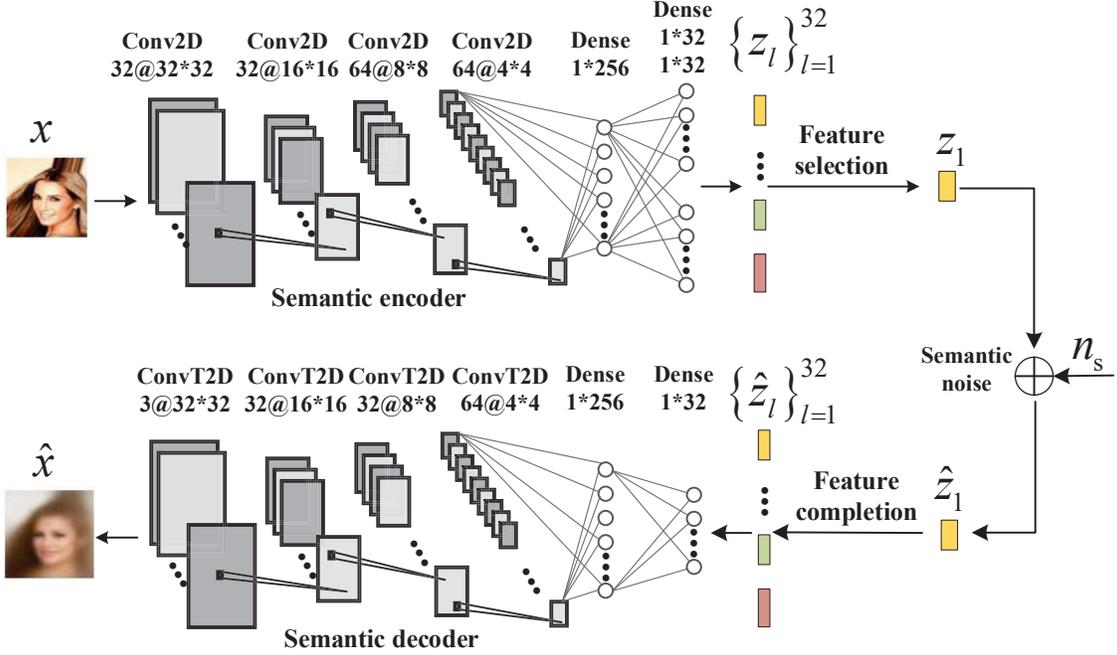


Fig. 4: Proposed robust  $\beta$ -VAE based architecture for the explainable semantic communication system

Note that, in the robust  $\beta$ -VAE network, we let  $\{\mu_l\}_{l=1}^L$  and  $\{\sigma_l\}_{l=1}^L$  denote the mean and the corresponding standard deviation of the approximate posterior  $q_\phi(z|x)$ , respectively. Moreover, a reparametrization trick [39] is applied to estimate gradients of the objective function (24) with respect to the parameter  $\phi$ , where random independent variables  $\{\varepsilon_l\}_{l=1}^L$  are sampled from a standard Gaussian distribution, i.e.,  $\varepsilon_l \sim \mathcal{N}(0, 1)$ . Then, the output features of the semantic encoder  $\{z_l\}_{l=1}^L$  are given as follows

$$z_l = \mu_l + \sigma_l \varepsilon_l, \quad l = 1, \dots, L. \quad (27)$$

Thus, the feature  $z_l$  is equivalent to being sampled from distribution  $\mathcal{N}(\mu_l, \sigma_l^2)$ , where  $l = 1, \dots, L$ .

### B. Feature Selection and Completion

With the disentangled and explainable features, the proposed semantic communications system further performs feature selection and completion at the transmitter and receiver, respectively. Specifically, since the receiver may only be interested in some of the features, the transmitter only sends the intended features  $\{z_l\}_{l \in \mathcal{L}_{\text{sel}}}$  according to their semantic meanings, rather than all of the extracted features  $\{z_l\}_{l \in \mathcal{L}}$ , which can further reduce the amount of information transmission.

For the receiver, the proposed semantic source and channel decoder include semantic feature completion and feature reconstruction. Specifically, for the unintended features subset are not transmitted  $\{z_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$ , the receiver generated the corresponding features  $\{\hat{z}_l\}_{l \in \mathcal{L} \setminus \mathcal{L}_{\text{sel}}}$  based on the receiver knowledge base, where both the dimensions and value ranges of sets  $z_l$  and  $\hat{z}_l$  are the same.

Then, according to the completed semantic features  $\hat{Z} = \{\hat{z}_l\}_{l \in \mathcal{L}}$ , the feature reconstruction module recovers the original data  $\hat{X}$ .

### C. Proposed Architecture

The proposed lightweight semantic communication architecture includes a semantic encoder network and a semantic decoder network, as shown in Fig. 4, where the notation Conv2D 32@32\*32 means that the network has 32 2-D convolutional filters of size 32\*32, and Dense 1\*256 represents a dense layer with 256 neurons. The details of the semantic encoder and the decoder network architectures are given as:

1) *Semantic encoder architecture*: : Conv2D 32@32\*32  $\rightarrow$  Conv2D 32@16\*16  $\rightarrow$  Conv2D 64@8\*8  $\rightarrow$  Conv2D 64@4\*4  $\rightarrow$  Dense 1\*256  $\rightarrow$  2 parallel Dense 1\*32  $\rightarrow \{z_l\}_{l=1}^{32} \rightarrow \{z_l\}_{l \in \mathcal{L}_{\text{sel}}}$ ;

2) *Semantic decoder architecture*: :  $\{\hat{z}_l\}_{l \in \mathcal{L}_{\text{sel}}} \rightarrow \{\hat{z}_l\}_{l=1}^{32} \rightarrow$  Dense 1\*32  $\rightarrow$  Dense 1\*256  $\rightarrow$  ConvT2D 64@4\*4  $\rightarrow$  ConvT2D 32@8\*8  $\rightarrow$  ConvT2D 32@16\*16  $\rightarrow$  ConvT2D 3@32\*32.

Note that, based on the feature selection, the proposed semantic communication system only needs to send the features  $\{z_l\}_{l \in \mathcal{L}_{\text{sel}}}$  that the receiver is interested in, instead of sending all features  $\{z_l\}_{l=1}^{32}$ .

## V. PROTOTYPE AND IMPLEMENTATIONS

The proposed architecture and hardware platform design of the semantic communication system prototype are shown in Fig. 5 (a) and (b), which can be used to implement the proposed robust and explainable semantic communications system in Fig. 3. The prototype includes two semantic

communication mobile users A and B. The trained robust  $\beta$ -VAE network is implemented at the portable RaspberryPi 4 Model B processors to realize the semantic encoding/decoding and the feature selection/completion functions. The integrated Wi-Fi module fulfills the bit-level transmission. The decoded data can be shown through the display.

The detailed parameters of the prototype are provided in Table III. The Raspberry Pi is installed with an ARM Cortex-A72@quad-core 1.5GHz CPU and 4GB of DDR4 RAM, and is equipped with Pytorch-CPU and torchvision software. The communication between Raspberry Pi A and B is realized through WiFi, where the socket is used to send and receive data, and Visdom is used to realize visual communication.

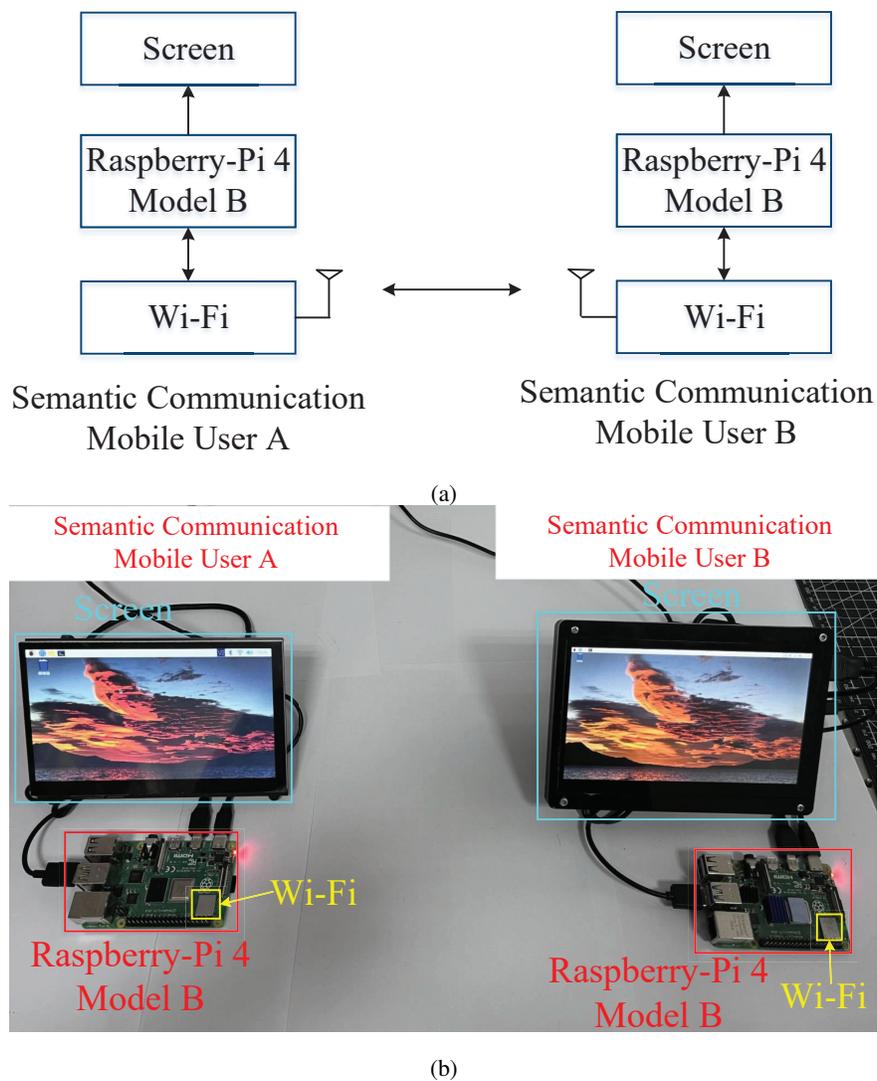


Fig. 5: (a) The architecture of the semantic communication system prototype; (b) The hardware platform of the semantic communication system prototype.

## VI. EXPERIMENTS AND DISCUSSIONS

In this section, we evaluate the proposed explainable semantic communications system using a graphics processing unit (GPU) and Raspberry Pi prototype, respectively. The GPU experiments in this work have been performed on 32 GB RAM i5-12600H, and 8 GB Nvidia GeForce 3060Ti GTX graphics card with Pytorch powered with CUDA 11.3. The experiments are performed via two standard datasets, i.e., MNIST Dataset and CelebA Dataset.

### A. Demonstration via GPU

First, we evaluate the robustness of the proposed semantic communication system. Specifically, the peak signal-to-noise ratio (PSNR) performance of the proposed robust  $\beta$ -VAE scheme with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  are demonstrated over the two channel models: the ANGC and a slow Rayleigh fading channel, where  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  mean that the trained SNRs of the schemes are 4dB and 8dB, respectively. Moreover, the PSNR performance of the deep joint source-channel coding (Deep-JSCC) scheme [19],  $\beta$ -VAE scheme, and the JPEG compression scheme are presented for comparisons.

Fig. 6 (a) shows PSNR versus different test SNRs of the four schemes over ANGC, where semantic noise parameters  $a = -0.1$ ,  $b = 0.1$  and  $\sigma_p^2 = 1$ . We observe that the PSNR of JPEG compression is the lowest among the five schemes, and the PSNR of the robust  $\beta$ -VAE schemes are higher than those of both Deep-JSCC and  $\beta$ -VAE. In the low SNR regions, the PSNR of the robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 4\text{dB}$  is the highest, and the PSNR of the robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 8\text{dB}$  is the higher than that of  $\beta$ -VAE, which verifies the robustness of our proposed design. Since the training noise of  $\text{SNR}_{\text{train}} = 4\text{dB}$  is higher than that of  $\text{SNR}_{\text{train}} = 8\text{dB}$ , the performance of  $\text{SNR}_{\text{train}} = 4\text{dB}$  is more robust, and thus the PSNR of  $\text{SNR}_{\text{train}} = 4\text{dB}$  is higher. In the high SNR regions, the PSNR of  $\beta$ -VAE, and robust  $\beta$ -VAE models tend to be the same. The reason is that the effect of noise at high SNR can be ignored.

TABLE III: Hardware parameters of the semantic communication prototype.

|                |                                     |
|----------------|-------------------------------------|
| GPU            | 500MHz VideoCore VI                 |
| CPU            | quad-core Cortex-A72                |
| System on Chip | Broadcom BCM2711@ 1.5GHz            |
| memory         | 4GB DDR4                            |
| Wi-Fi          | 2.4/ 5.0 GHz IEEE 802.11ac wireless |
| Screen         | 800 × 480 display                   |

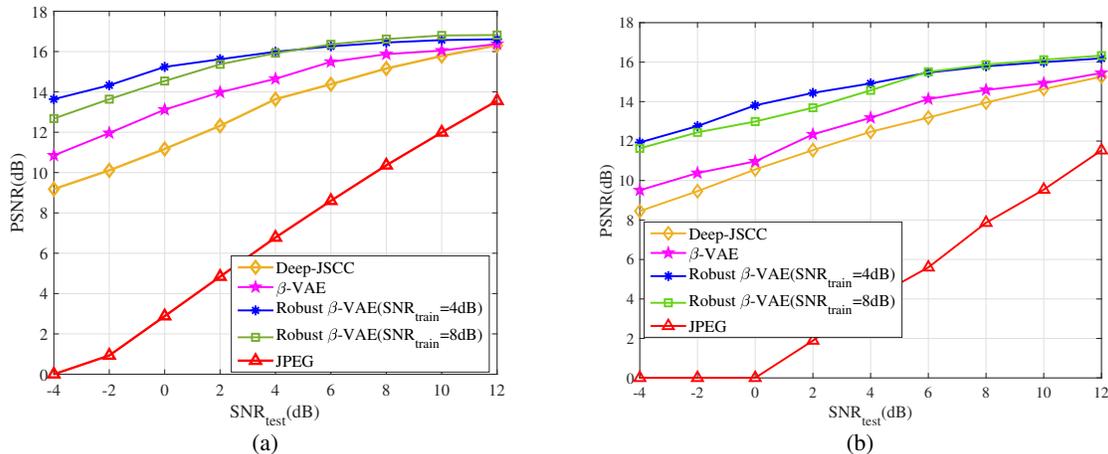


Fig. 6: (a) PSNR of Deep-JSCC,  $\beta$ -VAE, JPEG compression, and robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  over ANGC with semantic noise  $a = -0.1$ ,  $b = 0.1$  and  $\sigma_p^2 = 1$ ; (b) PSNR of Deep-JSCC,  $\beta$ -VAE, JPEG compression, and robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  over Rayleigh fading channels with  $\sigma_h^2 = 1$ ,  $a = -0.1$ ,  $b = 0.1$  and  $\sigma_p^2 = 1$ .

Fig. 6 (b) illustrates PSNR versus different test SNRs of the five schemes over the Rayleigh fading channel. Similar to Fig. 6 (a), the PSNR of JPEG compression is the lowest among the four schemes, and the PSNR of the robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  are higher than those of both Deep-JSCC and  $\beta$ -VAE. Note that for  $\text{SNR}_{\text{test}} = 8\text{dB}$ , the PSNR of the robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 8\text{dB}$  is the higher than that of  $\text{SNR}_{\text{train}} = 4\text{dB}$ . This is because the training SNR of  $\text{SNR}_{\text{train}} = 8\text{dB}$  is also  $8\text{dB}$ . Comparing Fig. 6 (a) with ANGC, the PSNRs of the schemes in Fig. 6 (b) are lower due to Rayleigh random fading.

Table IV illustrates the transmission performance of JPEG compression,  $\beta$ -VAE, and robust  $\beta$ -VAE with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$  over ANGC with semantic noise parameters  $a = -0.1$ ,  $b = 0.1$  and  $\sigma_p^2 = 1$ . The second column of Table IV shows the transmission performance of the JPEG compression scheme, where the transmitted semantics cannot be recognized from the received image. The third column shows the results of the  $\beta$ -VAE scheme, where the transmission semantics can be recognized from the received image. The fourth and fifth columns show received images of the robust  $\beta$ -VAE scheme with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$ , and the quality is better than that of the  $\beta$ -VAE scheme.

Table V illustrates transmission performance of the four schemes over the Rayleigh fading channel with semantic noise parameters  $a = -0.1$ ,  $b = 0.1$  and  $\sigma_p^2 = 1$ . The second column of table IV shows the transmission performance of the JPEG compression scheme, where the transmission semantics cannot be recognized from the received image. The third column shows

TABLE IV: Transmission performance comparison over ANG C

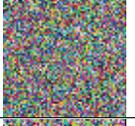
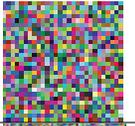
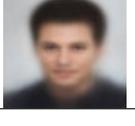
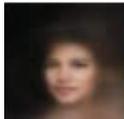
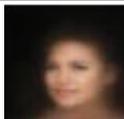
| Test SNR     | Data $X$  | JPEG  | $\beta$ -VAE  | Robust $\beta$ -VAE<br>$\text{SNR}_{\text{train}} = 4\text{dB}$                    | Robust $\beta$ -VAE<br>$\text{SNR}_{\text{train}} = 8\text{dB}$                     |
|--------------|---|---|---|--|---|
| SNR_test=4dB |  |  |  |  |  |
| SNR_test=8dB |  |  |  |  |  |
| SNR_test=4dB |  |  |  |  |  |
| SNR_test=8dB |  |  |  |  |  |

TABLE V: Transmission performance comparison over Rayleigh fading channel

| Test SNR     | Data $X$  | JPEG  | $\beta$ -VAE  | Robust $\beta$ -VAE<br>$\text{SNR}_{\text{train}} = 4\text{dB}$                      | Robust $\beta$ -VAE<br>$\text{SNR}_{\text{train}} = 8\text{dB}$                       |
|--------------|---|---|---|--|---|
| SNR_test=4dB |  |  |  |  |  |
| SNR_test=8dB |  |  |  |  |  |
| SNR_test=4dB |  |  |  |  |  |
| SNR_test=8dB |  |  |  |  |  |

the transmission performance of the  $\beta$ -VAE scheme, where the transmission semantics can be recognized from the received image. The fourth and fifth columns show the transmission performance of the robust  $\beta$ -VAE scheme with  $\text{SNR}_{\text{train}} = 4\text{dB}$  and  $\text{SNR}_{\text{train}} = 8\text{dB}$ , and the received image is better than that of the  $\beta$ -VAE scheme.

TABLE VI: Proposed semantic communication with feature selection

| Intended Feature                   | Skin color  | Face orientation  | Gender   | Hairstyle   |
|------------------------------------|---|---|--|---|
| Source data $X$                    |  |  |  |  |
| Receiver Knowledge Base            |  |  |  |  |
| Decoded data $\hat{X}$ (GPU)       |  |  |  |  |
| Decoded data $\hat{X}$ (Raspberry) |  |  |  |  |

### B. Demonstration via Prototype

In this subsection, we demonstrate that the proposed explainable semantic communication system with feature selection can improve the transmission efficiency via our prototype.

Table VI shows the performance of the proposed explainable semantic communication system with feature selection. From Column 2 to Column 5, we present four examples to show how the explainable encoder and feature selection work in the transmission. In the second column, the intended feature to send is skin color. The proposed semantic communication system performs feature extraction on the input white-skinned women picture, and then only selects the white skin color feature for transmission. Although the woman in the receiving knowledge base has darker skin, the reconstructed image is changed to white-skin. In the third column, the intended feature is face orientation. The proposed semantic communication system can successfully reconstruct a picture with the same face orientation at the receiver. Similarly, the intended features of the third and fourth columns are gender and hairstyle, respectively, and the proposed semantic communication system can also recover the correct feature at the receiver.

Table VII compares a compression ratio, transmission time, PSNR and reconstructed image of the original image transmission scheme, JPEG compression scheme,  $\beta$ -VAE scheme, and robust  $\beta$ -VAE scheme over our proposed semantic communication prototype on MNIST dataset with high SNR. From Table VII, we observe that the compression ratio of the  $\beta$ -VAE scheme and

robust  $\beta$ -VAE scheme is 78.4 which is significantly higher than those of the JPEG compression scheme (1.81) and original image transmission scheme. Thus, the transmission time of the  $\beta$ -VAE scheme and robust  $\beta$ -VAE scheme is about 0.3ms, which is significantly lower than those of the JPEG compression scheme (10.88ms) and original image transmission scheme (18.86ms). Therefore, the proposed semantic communication system can significantly reduce the transmission load and time. Moreover, the PSNR of the robust  $\beta$ -VAE scheme is close to that of the JPEG compression scheme, and is higher than that of the  $\beta$ -VAE scheme. Comparing of reconstructed images, we can clearly and accurately identify the number “7” from the recovered images using the proposed robust  $\beta$ -VAE scheme.

TABLE VII: Performance of the proposed semantic communication prototype on MNIST dataset

|                     | Transmission time (ms) | Compression ratio | PSNR  | Reconstructed image   |
|---------------------|------------------------|-------------------|-------|---|
| Original image      | 18.86                  | 1                 | 100   |    |
| JPEG                | 10.88                  | 1.81              | 17.49 |   |
| $\beta$ -VAE        | 0.30                   | 78.4              | 15.79 |  |
| Robust $\beta$ -VAE | 0.30                   | 78.4              | 16.07 |  |

Table VIII compares a compression ratio, transmission time, PSNR and reconstructed image of the original image transmission scheme, JPEG compression scheme,  $\beta$ -VAE scheme, and robust  $\beta$ -VAE scheme over our proposed semantic communication prototype on CelebA dataset with high SNR. Similar to Table VII, the compression ratio of the  $\beta$ -VAE scheme and robust  $\beta$ -VAE scheme is 384 which is significantly higher than those of the JPEG compression scheme (4.49) and original image transmission scheme. Thus, the transmission time of the  $\beta$ -VAE scheme and robust  $\beta$ -VAE scheme is about 0.18ms, which is significantly lower than those of the JPEG compression scheme (9.53ms) and original image transmission scheme (28.38ms). Therefore, the proposed semantic communication system can significantly reduce the transmission load and time. Moreover, the PSNR of the robust  $\beta$ -VAE scheme is close to that of the JPEG

TABLE VIII: Performance of the proposed semantic communication prototype on CelebA dataset

|                     | Transmission time (ms) | Compression ratio | PSNR  | Reconstructed image   |
|---------------------|------------------------|-------------------|-------|---|
| Original image      | 28.38                  | 1                 | 100   |  |
| JPEG                | 9.53                   | 4.49              | 30.17 |  |
| $\beta$ -VAE        | 0.18                   | 384               | 17.66 |  |
| Robust $\beta$ -VAE | 0.18                   | 384               | 19.73 |  |

compression scheme, and is higher than that of the  $\beta$ -VAE scheme. Note that, although the effect of the reconstructed image of proposed robust  $\beta$ -VAE scheme is a bit blurry, the three main semantic features of the original image: female, white skin color and long hair, are all accurately transmitted, which verifies the validity and accuracy of the proposed task-oriented semantic communication scheme.

## VII. CONCLUSIONS

In this paper, we propose an explainable and easy-to-implement semantic communication framework that is compatible with conventional communication systems. In this new framework, the semantic encoder can extract feature vectors, disentangle the semantic information, and improve robustness against semantic information ambiguity. To further reduce the communication cost, we apply feature selection to choose only task-related semantic information to transmit. Then, we present two information theoretic metrics, namely, the rate-distortion-perception function and semantic channel capacity to characterize the semantic information compression and transmission, respectively. To quantify the semantic information transmission with the additive quantization noise and physical channel noise, we further derive upper and lower bounds on the semantic channel capacity. Then, we propose a feasible design of the explainable semantic communication system, which includes a robust  $\beta$ -VAE lightweight unsupervised learning network. Finally, we develop a wireless mobile semantic communication

proof-of-concept prototype to implement the semantic communication design. Our experiments demonstrate that the proposed semantic communication system significantly outperforms the state-of-the-art methods, and shows robustness against various noise levels on two benchmark datasets. This work attempts to provide frameworks and theoretic metrics to explain and analyze the black-box semantic communications problem, and to provide guidelines on implementing the semantic communication in practical communication systems.

## VIII. APPENDICES

### APPENDIX A

#### PROOF OF LEMMA 1

We first derive the optimal conditional distribution  $q(\hat{x}|x)$  in (14) for a given output distribution  $r(x)$ . The mutual information  $I(X; \hat{X}) = \sum_x \sum_{\hat{x}} p(x) q(\hat{x}|x) \log \frac{q(\hat{x}|x)}{r(\hat{x})}$  is convex in  $q(\hat{x}|x)$  for fixed  $p(x)$ , and the KL divergence  $d_{KL}(p(x), r(\hat{x})) = \sum_x p(x) \log \frac{p(x)}{r(\hat{x})}$  is also convex in  $q(\hat{x}|x)$  for fixed  $p(x)$ . Thus, problem (14) is convex in  $q(\hat{x}|x)$ . Then, the Lagrangian function of problem (14) is given by

$$\begin{aligned} L(q(\hat{x}|x)) &= I(X; \hat{X}) + \alpha \sum_x \sum_{\hat{x}} p(x) q(\hat{x}|x) (x - \hat{x})^2 \\ &\quad + \mu \sum_x p(x) \log \frac{p(x)}{r(\hat{x})} + \sum_x \gamma(x) \sum_{\hat{x}} q(\hat{x}|x), \end{aligned} \quad (28)$$

where  $\alpha \geq 0$ ,  $\mu \geq 0$  and  $\gamma(x) \geq 0$  are Lagrange multipliers attached with constraints (14b), (14c) and (14d), respectively. For given  $r(\hat{x})$ , the derivative of (28) with respect to  $q(\hat{x}|x)$  is given as

$$\frac{\partial L(q(\hat{x}|x))}{\partial q(\hat{x}|x)} = p(x) \left( \log \frac{q(\hat{x}|x)}{r(\hat{x})} + \alpha(x - \hat{x})^2 - \mu \frac{p(x)}{r(\hat{x})} + \frac{\gamma(x)}{p(x)} \right). \quad (29)$$

Let  $\frac{\partial L(q(\hat{x}|x))}{\partial q(\hat{x}|x)} = 0$ , then we obtain the optimal  $q(\hat{x}|x)$  as

$$q^*(\hat{x}|x) = r(\hat{x}) \exp \left( \mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2 - \frac{\gamma(x)}{p(x)} \right) \quad (30a)$$

$$= \frac{r(\hat{x})}{\tilde{\gamma}(x)} \exp \left( \mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2 \right), \quad (30b)$$

where  $\tilde{\gamma}(x) \triangleq \exp \left( \frac{\gamma(x)}{p(x)} \right)$ .

Since  $\sum_{\hat{x}} q(\hat{x}|x) = 1$ , we have

$$\sum_{\hat{x}} \frac{r(\hat{x})}{\tilde{\gamma}(x)} \exp\left(\mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2\right) = 1. \quad (31)$$

Furthermore, we obtain

$$\tilde{\gamma}(x) = \sum_{\hat{x}} r(\hat{x}) \exp\left(\mu \frac{p(x)}{r(\hat{x})} - \alpha(x - \hat{x})^2\right). \quad (32)$$

Substituting (32) into (30b), we obtain the optimal  $q^*(\hat{x}|x)$  as given in Lemma 1.

From [37], we find that given a fixed conditional distribution  $q(\hat{x}|x)$ , the optimal output distribution  $r(\hat{x})$  is  $r^*(x) \triangleq \sum_x p(x)q(\hat{x}|x)$ . We rewrite the proof below.

$$I(X; Z) = \sum_{x, \hat{x}} p(x)q(\hat{x}|x) \log \frac{p(x)q(\hat{x}|x)}{p(x)r(\hat{x})} \quad (33)$$

$$- \sum_{x, \hat{x}} p(x)q(\hat{x}|x) \log \frac{p(x)q(\hat{x}|x)}{p(x)r^*(\hat{x})} \quad (34)$$

$$= \sum_{\hat{x}} r^*(x) \log \frac{r^*(\hat{x})}{r(\hat{x})} \geq 0, \quad (35)$$

where the last inequality holds because of the non-negative property of KL divergence.

## REFERENCES

- [1] L. Knud, "State of the IoT 2020: 12 billion IoT connections, surpassing non-IoT for the first time," <https://iot-analytics.com/state-of-the-iot-2020-12-billion-iot-connections-surpassing-non-iot/>, 2020.
- [2] J. Antoniou, "Quality of experience and emerging technologies: Considering features of 5G, IoT, cloud and AI," in *Quality of Experience and Learning in Information Systems*, pp. 1–8. Springer, 2021.
- [3] W. Saad, M. Bennis, and M. Chen, "A vision of 6G wireless systems: Applications, trends, technologies, and open research problems," *IEEE Netw.*, vol. 34, no. 3, pp. 134–142, Oct. 2020.
- [4] E. Calvanese Strinati, S. Barbarossa, J. L. Gonzalez-Jimenez, D. Ktenas, N. Cassiau, L. Maret, and C. Dehos, "6G: The next frontier: From holographic messaging to artificial intelligence using subterahertz and visible light communication," *IEEE Veh. Technol. Mag.*, vol. 14, no. 3, pp. 42–50, Oct. 2019.
- [5] B. Mao, F. Tang, Y. Kawamoto, and N. Kato, "AI models for green communications towards 6G," *IEEE Commun. Surveys Tuts.*, vol. 24, no. 1, pp. 210–247, Nov. 2022.
- [6] K. Niu, J. Dai, S. Yao, S. Wang, Z. Si, X. Qin, and P. Zhang, "Towards semantic communications: A paradigm shift," *arXiv preprint arXiv:2203.06692*, 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*, 2022.
- [8] M. Kountouris and N. Pappas, "Semantics-empowered communication for networked intelligent systems," *IEEE Commun. Mag.*, vol. 59, no. 6, pp. 96–102, Jan. 2021.

- [9] M. Sana and E. Calvanese Strinati, "Learning semantics: An opportunity for effective 6G communications," *arXiv preprint arXiv:2202.11958*, 2021.
- [10] Y. L. G. Shi Y. Xiao. and X. Xie, "From semantic communication to semantic-aware networking: Model, architecture, and open problems," *IEEE Commun. Mag.*, vol. 59, no. 8, pp. 44–50, Aug. 2021.
- [11] X. Luo, H.-H. Chen, and Q. Guo, "Semantic communications: Overview, open issues, and future research directions," *IEEE Wirel. Commun.*, pp. 1–10, Jan. 2022.
- [12] J. Bao, P. Basu, M. Dean, C. Partridge, A. Swami, W. Leland, and J. A. Hendler, "Towards a theory of semantic communication," in *Proc. IEEE Netw. Sci. Workshop*, pp. 110–117, Jun. 2011.
- [13] A. Y. B. Güler and A. Swami, "The semantic communication game," *IEEE Trans. Cogn. Commun. Netw.*, vol. 4, no. 4, pp. 787–802, Dec. 2018.
- [14] W. Weaver, "Recent contributions to the mathematical theory of communication," *ETC: a review of general semantics*, pp. 261–281, Sep. 1949.
- [15] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, no. 3, pp. 623–656, Jul. 1948.
- [16] N. Farsad, M. Rao, and A. Goldsmith, "Deep learning for joint source-channel coding of text," in *Proc.(ICASSP)*, pp. 2326–2330, Apr. 2018.
- [17] H. Xie, Z. Qin, L. Geoffrey Ye., and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE Trans. Signal Process.*, vol. 69, pp. 2663–2675, Apr. 2021.
- [18] Q. Hu, G. Zhang, Z. Qin, Y. Cai, and G. Yu, "Robust semantic communications against semantic noise," *arXiv preprint arXiv:2202.03338*, Feb. 2022.
- [19] E. Bourtsoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Trans. Cognit. Commun. Netw.*, vol. 5, no. 3, pp. 567–579, May. 2019.
- [20] D. B. Kurka and D. Gündüz, "Deepjpsc-f: Deep joint source-channel coding of images with feedback," *IEEE J. Sel. Areas Inf. Theory*, vol. 1, no. 1, pp. 178–193, Apr. 2020.
- [21] M. Yang, C. Bian, and H.-S. Kim, "OFDM-guided deep joint source channel coding for wireless multipath fading channels," *IEEE Trans. Cognit. Commun. Netw.*, Feb. 2022.
- [22] K. Choi, K. Tatwawadi, A. Grover, T. Weissman, and S. Ermon, "Neural joint source-channel coding," in *Proc. Int. Conf. Mach. Learn.(ICML)*. PMLR, pp. 1182–1192, Jun. 2019.
- [23] 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 Trans. Circuits Syst. Video Technol.*, May. 2021.
- [24] H. Tong, Z. Yang, S. Wang, Y. Hu, W. Saad, and C. Yin, "Federated learning based audio semantic communication over wireless networks," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, pp. 1–6, Feb. 2021.
- [25] Z. Weng and Z. Qin, "Semantic communication systems for speech transmission," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2434–2444, Aug. 2021.
- [26] H. Xie, Z. Qin, and G. Y. Li, "Task-oriented multi-user semantic communications for VQA," *IEEE Wirel. Commun. Lett.*, vol. 11, no. 3, pp. 553–557, Dec. 2021.
- [27] Q. Zhou, R. Li, Z. Zhao, Y. Xiao, and H. Zhang, "Adaptive bit rate control in semantic communication with incremental knowledge-based HARQ," *arXiv preprint arXiv:2203.06634*, 2022.
- [28] P. Jiang, C.-K. Wen, S. Jin, and G. Y. Li, "Deep source-channel coding for sentence semantic transmission with HARQ," *arXiv preprint arXiv:2106.03009*, 2021.
- [29] Z. Z. C. P. Q. Zhou R. Li. and H. Zhang, "Semantic communication with adaptive universal transformer," *IEEE Wirel. Commun. Lett.*, vol. 11, no. 3, pp. 453–457, Dec. 2021.

- [30] H. Xie and Z. Qin, "A lite distributed semantic communication system for internet of things," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 142–153, Jan. 2021.
- [31] K. Lu, R. Li, X. Chen, Z. Zhao, and H. Zhang, "Reinforcement learning-powered semantic communication via semantic similarity," *arXiv preprint arXiv:2108.12121*, 2021.
- [32] M. Yang and H.-S. Kim, "Deep joint source-channel coding for wireless image transmission with adaptive rate control," *arXiv preprint arXiv:2110.04456*, 2021.
- [33] M. Ding, J. Li, M. Ma, and X. Fan, "SNR-adaptive deep joint source-channel coding for wireless image transmission," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.(ICASSP)*, pp. 1555–1559, May. 2021.
- [34] Y. M. J. Shao and J. Zhang, "Learning task-oriented communication for edge inference: An information bottleneck approach," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 1, pp. 197–211, Jan. 2022.
- [35] Y. Blau and T. Michaeli, "Rethinking lossy compression: The rate-distortion-perception tradeoff," in *International Conference on Machine Learning*. PMLR, pp. 675–685, 2019.
- [36] G. Zhang, J. Qian, J. Chen, and A. Khisti, "Universal rate-Distortion-Perception representations for lossy compression," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [37] T. M. Cover and J. A. Thomas, *Elements of information theory, 2nd ed.*, New York, NY, USA: Wiley, 2006.
- [38] S. Ihara, "On the capacity of channels with additive non-Gaussian noise," *Inform. Contr.*, vol. 37, no. 1, pp. 34–39, Sep. 1978.
- [39] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," *Proc. ICLR*, pp. 1–12, 2017.
- [40] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [41] S. Eguchi and S. Kato, "Entropy and divergence associated with power function and the statistical application," *Entropy*, vol. 2, pp. 262–274, Dec. 2010.
- [42] F. Futami, I. Sato, and M. Sugiyama, "Variational inference based on robust divergences," *arXiv preprint arXiv:1710.06595*, 2017.