# Adaptive Information Bottleneck Guided Joint Source and Channel Coding for Image Transmission

Lunan Sun, Yang Yang, *Member, IEEE*, Mingzhe Chen, *Member, IEEE*, Caili Guo, *Senior Member, IEEE*, Walid Saad, *Fellow, IEEE* and H. Vincent Poor, *Life Fellow, IEEE*

*Abstract*—Joint source and channel coding (JSCC) for image transmission has attracted increasing attention due to its robustness and high efficiency. However, the existing deep JSCC research mainly focuses on minimizing the distortion between the transmitted and received information under a fixed number of available channels. Therefore, the transmitted rate may be far more than its required minimum value. In this paper, an adaptive information bottleneck (IB) guided joint source and channel coding (AIB-JSCC) method is proposed for image transmission. The goal of AIB-JSCC is to reduce the transmission rate while improving the image reconstruction quality. In particular, a new IB objective for image transmission is proposed so as to minimize the distortion and the transmission rate. A mathematically tractable lower bound on the proposed objective is derived, and then, adopted as the loss function of AIB-JSCC. To trade off compression and reconstruction quality, an adaptive algorithm is proposed to adjust the hyperparameter of the proposed loss function dynamically according to the distortion during the training. Experimental results show that AIB-JSCC can significantly reduce the required amount of transmitted data and improve the reconstruction quality and downstream task accuracy.

*Index Terms*—Information bottleneck, joint source and channel coding, image transmission.

## I. INTRODUCTION

Shannon's information theory has laid the foundations of modern communication systems. In particular, according to Shannon's information theory, separate source and channel coding (SSCC) is optimal for a memoryless source and channel when the latency, complexity, and code length are not constrained [1]. However, SSCC has several practical limitations. First, the theory is based on the assumption of potentially infinite code lengths, which are impossible in practice, and SSCC is suboptimal for finite code lengths. Also, to achieve theoretically optimal performance, maximum likelihood detection methods must be used, which can be, in general, NP-hard [2], thus introducing very high computational complexity and leading to unacceptable latency. Furthermore, the envisioned sixth generation (6G) of wireless networks are expected to connect trillion-level devices and require 10 to 1000 times higher rates [3]. In addition, it is thought that 6G will support a wide range of services and applications [4], such as ugmented reality, medical imaging and autonomous vehicles [5], [6], which have strict latency requirement [7]–[10]. Therefore, SSCC may not be able to meet the requirements of 6G.

To address the above-mentioned challenges, joint source and channel coding (JSCC) has attracted increasing attention as a means to achieve reliable data transmission. Existing studies of JSCC can be classified into two types: traditional research based on mathematical models [11]–[14] and deep learning (DL)-based research [15]–[19]. Traditional JSCC research mainly relies on traditional source coding and channel coding theory while focusing on performance analysis under ideal channel or source assumptions [11], [12]. Coding schemes, such as bit allocation algorithm [13], robust nonlinear block coding [14] have also been studied. However, these hand-crafted coding schemes may require additional tuning. Motivated by the impressive performance of DL in many domains such as computer vision [20], image compression [21], and natural language processing [22], DL-based JSCC has been extensively studied [15]–[19], which can potentially support future semantic communications [23], [24]. Specifically, since images have larger dimensions than speech and text data, there is more information redundancy in images, and transmitting image data requires higher rate than transmitting speech and text data. Therefore, it is more challenging to design a DL-based JSCC system for image transmission.

### A. Related Works and Challenges

The existing works on DL-based JSCC for image transmission model the communication system as a deep neural network (DNN)-based autoencoder [15]–[19]. The main goal is jointly training the encoder and decoder to preserve information and improve the reconstruction quality. Minimizing the mean-squared error (MSE) between the input images

L. Sun and Y. Yang are with the Beijing Key Laboratory of Network System Architecture and Convergence, School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: sunlunan@bupt.edu.cn; young0607@bupt.edu.cn).

M. Chen is with the Department of Electrical and Computer Engineering and Institute for Data Science and Computing, University of Miami, Coral Gables, FL, 33146 USA (e-mail: mingzhe.chen@miami.edu).

C. Guo is with the Beijing Laboratory of Advanced Information Net works, School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: guocaili@bupt.edu.cn).

W. Saad is with the Wireless@VT Group, Bradley Department of Electrical and Computer Engineering, Virginia Tech, Arlington, VA, USA (email: walids@vt.edu).

H. V. Poor is with the Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: poor@princeton.edu)

arXiv:2203.06492v2 [cs.IT] 29 May 2023

and output reconstructions [15]–[17] is commonly used to achieve this goal. In particular, the authors in [15] proposed an autoencoder-based JSCC architecture called deep JSCC that minimizes the MSE between the original images and the recovered images. Deep JSCC outperforms SSCC that combines JPEG or JPEG2000 with capacity-achieving channel codes. In [16], the authors incorporated the channel output feedback into the transmission system and further improved the reconstruction quality of Deep JSCC. To address the variations of signal-to-noise ratios (SNRs) during transmission, the work in [17] designed a novel JSCC scheme, which uses a channel-wise soft attention network to adapt automatically to various channel conditions. These existing works [15]–[17] that use MSE as the distortion function to recover each pixel equally in image transmission may lose the information of important pixels thus reducing image reconstruction quality. In contrast, mutual information measures the distortion in terms of the distribution of images, which can emphasize key pixels and has stronger generalization ability. In [18], a discrete variational autoencoder model is designed to maximize the mutual information between the source and noisy codewords. The authors in [19] developed a JSCC model to maximize the mutual information between the codewords and input image. Overall, the existing works on DL-based JSCC for image transmission aim at minimizing the distortion between the transmitted and received images by utilizing various distortion metrics such as MSE and mutual information as loss functions under a fixed number of achievable channels. While the works in [15]–[19] are interesting, the theoretical minimum description length (or transmission rate) of the codewords to express source is neglected in the loss function. Therefore, the transmission rate may be much larger than the minimum required rate. A new form of loss function for JSCC, that simultaneously minimizes the transmission rate and the distortion deserves investigation.

Recently, the authors in [25] proposed an information-theoretic principle, termed information bottleneck (IB) to compress information and improve data fitting performance simultaneously by using mutual information between the codewords and the labels of the inputs as distortion. IB principle has been extensively applied in many domains including improving the performance of generalization and robustness [26], suppressing irrelevant features [27], and dealing with domain shift [28]. Since IB inherits the properties of RD theory, it can characterize the maximal compression ratio and the optimal features in theory [29], [30]. Therefore, we propose a novel IB-guided JSCC that can reduce the transmission rate for a given reconstruction quality. Here, we need to note that it is challenging to apply the IB principle in JSCC for image transmission since standard IB is particularly designed for supervised tasks, while a JSCC-based image transmission system can be viewed as an unsupervised data reconstruction task. Meanwhile, in an image transmission JSCC system, the distribution of the input images is usually unknown, and the dimension of the extracted codewords is large. Thus, the mutual information used in IB is intractable. Therefore, to apply the IB principle to JSCC for image transmission, two main challenges must be addressed:

- How to design a proper form of IB for an image transmission JSCC system, which is unsupervised.
- How to calculate the mutual information used in IB and obtain a tractable and differentiable IB objective.

### B. Contributions

The main contribution of this paper is an adaptive IB-guided JSCC (AIB-JSCC) scheme for image transmission to address the above issues. The major contributions of the paper can be summarized as follows:

- We design a new form of IB objective that aims at simultaneously maximizing the mutual information between the received noisy codewords and the input images, and minimizing the mutual information between the transmitted codewords and the input images. Thus, the new IB objective enables the image transmission JSCC system to reduce the transmission rate while guaranteeing the reconstruction quality. To the best of the authors' knowledge, this is the first work that applies the IB principle to image transmission JSCC and provides a theoretically maximal compression ratio guidance for neural networks.
- As the mutual information in the proposed IB objective is intractable for DNNs with high-dimensional features, we develop a new mathematically tractable and differentiable lower bound on the proposed IB objective via a variational lower bound and contrastive log-ratio upper bound (CLUB) on mutual information. The derived lower bound is used as the loss function of AIB-JSCC.
- We propose an adaptive algorithm, which can adjust the hyperparameter of the proposed IB objective to balance the reconstruction distortion and the required transmission rate. In particular, we first develop an algorithm to adjust the hyperparameter value dynamically by exploiting reconstruction error. Then, we derive an upper bound on the hyperparameter, which can prevent excessive information discarding in the transmitted codewords.

We compare AIB-JSCC with traditional SSCC and state-of-the-art JSCC methods and quantify the performance gain via extensive experiments. Simulation results show that AIB-JSCC significantly reduces the required storage space and the amount of transmitted image data.

The rest of this paper is organized as follows. In Section II, the system model is described. The proposed IB objective is presented in Section III. The adaptive IB algorithm is introduced in Section IV. In Section V, we provide extensive experimental results to verify the effectiveness of AIB-JSCC. Finally, the conclusions are drawn in Section VI.

## II. SYSTEM MODEL

In this section, we first describe the studied JSCC system model for image transmission. Then, we discuss the motivation for our work as well as the IB principle.

### A. System Model

As shown in Fig. 1, we consider a point-to-point image transmission system [15]–[19]. An input image with size $H$

TABLE I
LIST OF NOTATION.

| Notation | Definition | Notation | Definition |
|---|---|---|---|
| $N$ | The size of the images | $M$ | The length of codewords |
| $B$ | The sample number in a batch | $P$ | The number of parallel channels |
| $\text{MSE}[w]$ | The MSE between the inputs and the reconstructions at the $w$-th epoch | $\eta_\varepsilon(\cdot)$ | Th transition function of BSC with error probability $\varepsilon$ |
| $\varphi$ | The parameters of the encoder neural network | $\theta$ | The parameters of the decoder neural network |
| $\varepsilon$ | The error probability of the channel | $\varepsilon_k$ | The error probability of the $k$-th subchannel |
| $\boldsymbol{x}$ | The input images | $\boldsymbol{x}^{(i)}$ | The $i$-th input image |
| $\hat{\boldsymbol{x}}$ | The recovered images | $\hat{\boldsymbol{x}}^{(i)}$ | The $i$-th recovered image |
| $x_j^{(i)}$ | The $j$-th pixel in the $i$-th input image | $\hat{x}_j^{(i)}$ | The $j$-th pixel in the $i$-th recovered image |
| $\boldsymbol{y}$ | The codewords to be transmitted | $\boldsymbol{y}^{(i)}$ | The codewords extracted from $\boldsymbol{x}^{(i)}$ |
| $y_m$ | The $m$-th element in $\boldsymbol{y}$ | $y_m^{(i)}$ | The $m$-th element in $\boldsymbol{y}^{(i)}$ |
| $\hat{\boldsymbol{y}}$ | The noisy codewords received by decoder | $\hat{\boldsymbol{y}}^{(i)}$ | The noisy codewords extracted from $\boldsymbol{x}^{(i)}$ |
| $\hat{y}_m$ | The $m$-th element in $\hat{\boldsymbol{y}}$ | $\hat{y}_m^{(i)}$ | The $m$-th element in $\hat{\boldsymbol{y}}^{(i)}$ |
| $\boldsymbol{y}_{\text{ch}k}$ | The subcodewords to be transmitted across the $k$-th subchannel | $\hat{\boldsymbol{y}}_{\text{ch}k}$ | The noisy subcodewords received across the $k$-th subchannel |
| $f_\varphi(\cdot)$ | The output of the encoder neural network | $E_\varphi(\cdot)$ | The encoding process from $\boldsymbol{x}$ to $\boldsymbol{y}$ |
| $g_\theta(\cdot)$ | The output of the decoder neural network | $D_\theta(\cdot)$ | The decoding process from $\hat{\boldsymbol{y}}$ to $\hat{\boldsymbol{x}}$ |
| $\beta$ | The hyperparameter of the proposed IB objective | $\beta_{\text{adp}}$ | The value of $\beta$ calculated by the adaptive IB algorithm |
| $\beta_{\max}$ | The upper bound on $\beta$ | $\beta_{\min}$ | The minimum value of $\beta$ |

(Height) $\times W$ (Width) $\times C$ (Channel) is represented as a vector $\boldsymbol{x} \in \mathbb{R}^N$, where $\mathbb{R}$ represents the set of real numbers and $N = H \times W \times C$. The encoder encodes the image $\boldsymbol{x}$ into a binary codeword $\boldsymbol{y} \in \{0, 1\}^M$, where $M$ represents the length of the codeword $\boldsymbol{y}$ to be transmitted. The encoding function $E_\varphi : \mathbb{R}^N \to \{0, 1\}^M$ is parameterized by an encoder neural network with parameters $\varphi$, and the encoding process can be expressed as

$$\boldsymbol{y} = E_\varphi(\boldsymbol{x}), \tag{1}$$

where $\boldsymbol{y}$ is the JSCC codeword generated by encoding the source information and adding redundancy for error protection jointly. $\boldsymbol{y}$ is then transmitted across a noisy channel. To simplify the analysis, we do not consider concrete modulation, detection and decision schemes, and only consider the transmission of the codeword through a channel with a certain error probability, i.e., memoryless binary symmetric channel (BSC)[1]. Here, we consider a BSC with error probability $0 \leq \varepsilon \leq 0.5$, denoted by $\eta_\varepsilon : \{0, 1\}^M \to \{0, 1\}^M$. The channel output noisy codeword $\hat{\boldsymbol{y}} \in \{0, 1\}^M$ received by the decoder is expressed as

$$\hat{\boldsymbol{y}} = \eta_\varepsilon(\boldsymbol{y}) = \boldsymbol{y} \oplus \boldsymbol{z}, \tag{2}$$

where $\boldsymbol{z} \sim \text{Bern}(\varepsilon)$ represents the Bernoulli distributed noise of the considered channel, and $\oplus$ represents modulo-

[1]The BSC is a well-established and widely-used model in communication theory and information theory [1], [31], [32] since it is a simple and tractable model for theoretical analysis.
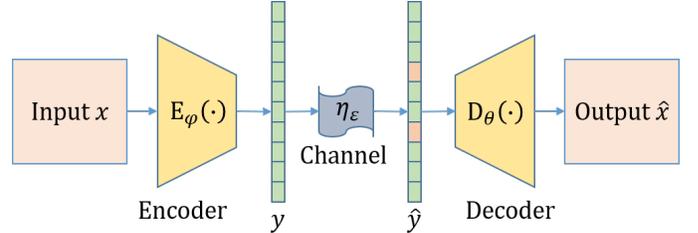


Fig. 1. An illustration of the JSCC system.

2 addition [33]–[35]. The channel capacity of the BSC with error probability $\varepsilon$ is

$$C_{\text{BSC}}(\varepsilon) = 1 - h(\varepsilon), \tag{3}$$

where $h(\varepsilon) = -\varepsilon \log \varepsilon - (1 - \varepsilon) \log(1 - \varepsilon)$ is the binary entropy function, and $\log(x) \triangleq \log_2(x)$.

The decoder decodes the noisy codeword $\hat{\boldsymbol{y}}$ into reconstructed image $\hat{\boldsymbol{x}} \in \mathbb{R}^N$. The decoding function is parameterized by the decoder neural network parameters $\theta$, and the decoding process is expressed as $D_\theta : \{0, 1\}^M \to \mathbb{R}^N$. The reconstructed image $\hat{\boldsymbol{x}}$ is

$$\hat{\boldsymbol{x}} = D_\theta(\hat{\boldsymbol{y}}) = D_\theta(\eta_\varepsilon(E_\varphi(\boldsymbol{x}))). \tag{4}$$

The goal of the considered system is to determine the encoder and decoder parameters that minimize the average reconstruction error between $\boldsymbol{x}$ and $\hat{\boldsymbol{x}}$ while keeping the
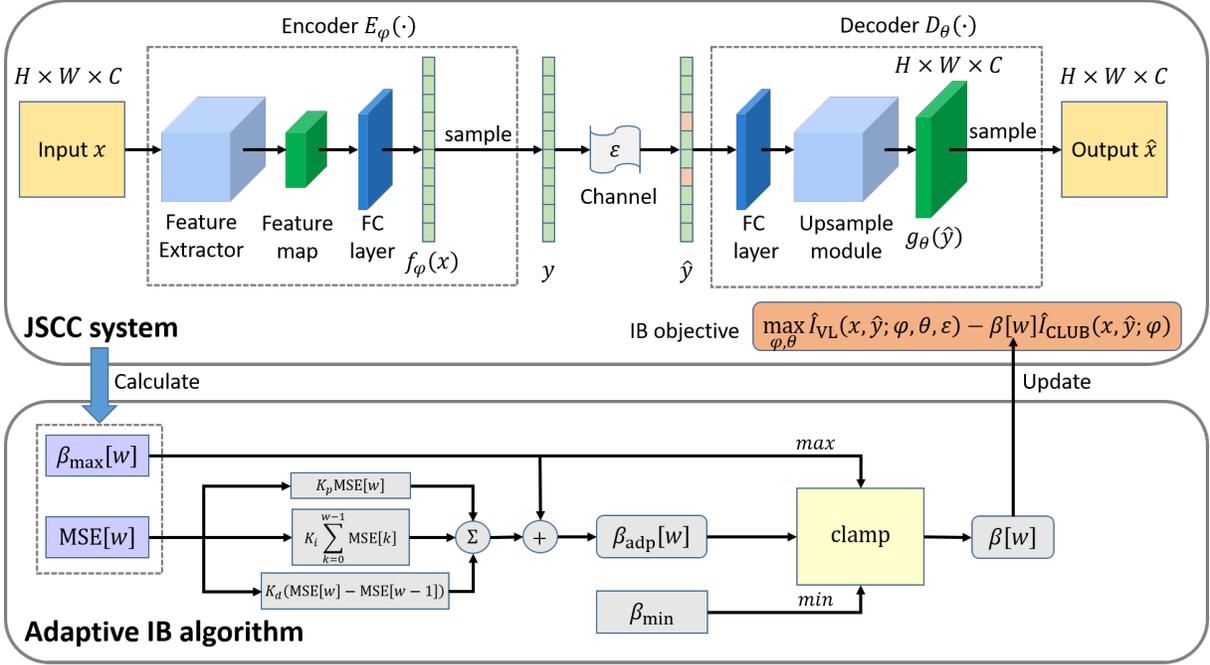
Fig. 2. An illustration of our proposed AIB-JSCC. Top: JSCC system with the proposed IB objective. Bottom: adaptive IB algorithm. First, we train the JSCC system which consists of an autoencoder, by optimizing the proposed IB objective. Then, we adjust $\beta$ according to the proposed algorithm. Finally, we alternately change $\beta$ and train the network.

minimum description length (or transmission rate) of $\boldsymbol{y}$ to express $\boldsymbol{x}$ short.

### B. Motivation

Existing deep JSCC solutions for image transmission [15]–[19] aim to minimize the distortion given a fixed number of available channels. However, they ignore the minimum description length (or transmission rate) of $\boldsymbol{y}$ to express $\boldsymbol{x}$, i.e., $I(\boldsymbol{x}; \boldsymbol{y})$ in the loss function. This motivates us to design a new loss function that can optimize both the distortion $d(\boldsymbol{x}, \hat{\boldsymbol{x}})$ and the transmission rate $I(\boldsymbol{x}; \boldsymbol{y})$ simultaneously. We resort to the IB principle. To extract the contained information of a target random variable $\boldsymbol{t}$ (e.g. label) in input $\boldsymbol{x}$, the authors in [25] used the mutual information between $\boldsymbol{y}$ and $\boldsymbol{t}$, $I(\boldsymbol{y}; \boldsymbol{t})$ as distortion measurement and proposed IB principle. The objective of IB is

$$\max_{p(\boldsymbol{y}|\boldsymbol{x})} \left[ I(\boldsymbol{y}; \boldsymbol{t}) - \beta I(\boldsymbol{x}; \boldsymbol{y}) \right]. \quad (5)$$

The first term $I(\boldsymbol{y}; \boldsymbol{t})$ in (5) encourages $\boldsymbol{y}$ to predict $\boldsymbol{t}$, and the second term $I(\boldsymbol{x}; \boldsymbol{y})$ in (5) encourages $\boldsymbol{y}$ to compress the information related to $\boldsymbol{x}$. According to (5), the system can obtain the optimal $\boldsymbol{y}$ that is maximally compressed with a certain distortion [30]. (5) is typically used as the loss function of supervised artificial intelligence tasks [26], [27], where $\boldsymbol{x}$ is the input image, $\boldsymbol{y}$ is the codeword, and $\boldsymbol{t}$ is the label of $\boldsymbol{x}$.

Even though the IB principle provides a new form of mutual information distortion, and can be used to guide the generation of the optimal features, the IB form in (5) is designed for supervised learning, and a label variable $\boldsymbol{t}$ is needed. Therefore, (5) cannot be applied to JSCC for image transmission directly, since image transmission is an unsupervised task. Moreover, the value of hyperparameter $\beta$ in (5) needs to be carefully designed to balance prediction and compression. To solve these problems, we propose a new form of the IB objective that can minimize both the distortion and transmission rate for image transmission JSCC. We derive a tractable and differentiable lower bound on the proposed objective and use the bound as the loss function of AIB-JSCC for image transmission. An adaptive algorithm is also designed to dynamically adjust the hyperparameter $\beta$, so as to balance the compression and reconstruction quality.

## III. PROPOSED IB OBJECTIVE FOR JSCC SYSTEM

This section first introduces the proposed IB objective for image transmission JSCC system. To obtain a tractable and differentiable form of the proposed IB objective, we then derive the lower bound of the proposed IB objective according to the variational lower bound and the upper bound of the mutual information.

### A. Proposed IB Objective

The considered JSCC system is shown in Fig. 2, and it mainly consists of an encoder $\mathrm{E}_{\boldsymbol{\varphi}}(\cdot)$ block, a decoder block $\mathrm{D}_{\boldsymbol{\theta}}(\cdot)$, a channel block, and an adaptive IB algorithm block. In the considered system, IB principle is adopted to guide JSCC to achieve theoretically minimal transmission rate with a certain tasks distortion. However, the standard form of the IB principle shown as (5) is not applicable to the considered system, since image transmission is an unsupervised task without label. To overcome this problem, we propose a new
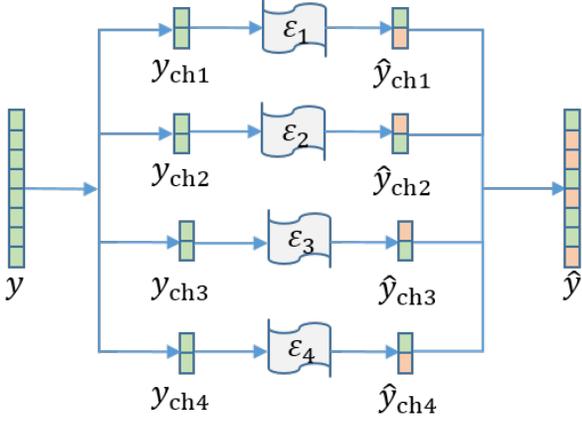
Fig. 3. An illustration of the parallel-channel case with 4 subchannels. The codeword $\boldsymbol{y}$ is first equally divided into 4 subcodewords, $\boldsymbol{y}_{\mathrm{ch1}}$, $\boldsymbol{y}_{\mathrm{ch2}}$, $\boldsymbol{y}_{\mathrm{ch3}}$ and $\boldsymbol{y}_{\mathrm{ch4}}$. These subcodewords are transmitted through their corresponding subchannel. At the receiver, the noisy subcodewords $\hat{\boldsymbol{y}}_{\mathrm{ch1}}$, $\hat{\boldsymbol{y}}_{\mathrm{ch2}}$, $\hat{\boldsymbol{y}}_{\mathrm{ch3}}$ and $\hat{\boldsymbol{y}}_{\mathrm{ch4}}$ are concatenated in order to obtain the noisy codeword $\hat{\boldsymbol{y}}$.

form of the IB objective for image transmission JSCC as follows:

$$\max_{\boldsymbol{\varphi},\boldsymbol{\theta}} \left[ I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right) - \beta I\left(\boldsymbol{x};\boldsymbol{y}\right) \right]. \quad (6)$$

In the proposed IB objective, we use $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ to capture the reconstruction distortion between $\boldsymbol{x}$ and $\hat{\boldsymbol{x}}$. By maximizing $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$, we can ensure that $\hat{\boldsymbol{y}}$ can capture the most useful information from $\boldsymbol{x}$. Hence, we can maximize $I\left(\boldsymbol{x};\hat{\boldsymbol{x}}\right)$, which is intractable due to unknown conditional probability $p\left(\boldsymbol{x}|\hat{\boldsymbol{x}}\right)$, via maximizing $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ since $\hat{\boldsymbol{x}}$ is reconstructed from $\hat{\boldsymbol{y}}$. However, since $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right) \leq I\left(\boldsymbol{x};\boldsymbol{y}\right)$, solely maximizing $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ may result in severe information redundancy in $\boldsymbol{y}$, which implies that the system requires much higher transmission rate to transmit $\boldsymbol{y}$. Thus, we use the second term, which is the transmission rate over the channel, to compress the information in $\boldsymbol{y}$. We minimize $I\left(\boldsymbol{x};\boldsymbol{y}\right)$ so that the minimum description length (or transmission rate) of $\boldsymbol{y}$ to express $\boldsymbol{x}$ can be reduced. Although the sizes of $\boldsymbol{x}$ and $\boldsymbol{y}$ are fixed, the probability distribution of $\boldsymbol{y}$ can be optimized to minimize the minimum description length (or transmission rate) of $\boldsymbol{y}$ that is used to represent $\boldsymbol{x}$. We treat the proposed loss function as a joint optimization problem that integrates both reconstruction distortion minimization and transmission rate minimization. Utilizing (6) as the loss function, the JSCC system can accurately transmit images while minimizing the required transmission rate.

However, (6) still cannot be applied to the JSCC systems, since the mutual information terms $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ and $I\left(\boldsymbol{x};\boldsymbol{y}\right)$ in (6) are mathematically intractable due to the unknown $p\left(\boldsymbol{x},\boldsymbol{y}\right)$, $p\left(\boldsymbol{x},\hat{\boldsymbol{y}}\right)$, $p\left(\boldsymbol{x}\right)$, $p\left(\boldsymbol{y}\right)$ and $p\left(\hat{\boldsymbol{y}}\right)$. To circumvent this challenge, we next derive the variational lower bound on $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ and estimate the upper bound on $I\left(\boldsymbol{x};\boldsymbol{y}\right)$.

*B. Variational Lower Bound on $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$*

Instead of maximizing the true value of $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$, we maximize its lower bound. We utilize the variational lower bound

on $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$, which is obtained by [36]

$$I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right) = \underbrace{H\left(\boldsymbol{x}\right)}_{\text{constant}} + \underbrace{\mathbb{E}_{p(\boldsymbol{x},\hat{\boldsymbol{y}})} \log \left[ \frac{p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)}{q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)} \right]}_{D_{KL}(p(\boldsymbol{x}|\hat{\boldsymbol{y}})||q(\boldsymbol{x}|\hat{\boldsymbol{y}})) \geq 0} \quad (7)$$
$$+ \underbrace{\mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x})} \mathbb{E}_{\hat{\boldsymbol{y}} \sim p(\hat{\boldsymbol{y}}|\boldsymbol{x})} \log \left[ q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right) \right]}_{I_{\mathrm{VL}}(\boldsymbol{x};\hat{\boldsymbol{y}})}.$$

In (7), $q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$ is the variational approximation of the true posterior $p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$. The first term, $H\left(\boldsymbol{x}\right)$ is the entropy of the input images, which is a constant and cannot be optimized by neural networks. The second term is the Kullback-Leibler (KL) divergence between $p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$ and $q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$, which is positive. Since $H\left(\boldsymbol{x}\right) \geq 0$ and $D_{KL}\left(p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)||q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)\right) \geq 0$, the variational lower bound on $I\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ is $I_{\mathrm{VL}}\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$, as defined in (7). The approximation error will be smaller if $q\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$ is closer to $p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$,

Since the conditional probability of $\hat{\boldsymbol{y}}$ given $\boldsymbol{x}$ depends on $\boldsymbol{\varphi}$ and $\boldsymbol{\varepsilon}$, we represent it as $p\left(\hat{\boldsymbol{y}}|\boldsymbol{x};\boldsymbol{\varphi},\boldsymbol{\varepsilon}\right)$. Denote the conditional probability of $\hat{\boldsymbol{x}}$ given $\hat{\boldsymbol{y}}$ as $p_{\boldsymbol{\theta}}\left(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}\right)$, which is parameterized by the decoder neural network. We use $p_{\boldsymbol{\theta}}\left(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}\right)$ as the variational approximation of the true posterior $p\left(\boldsymbol{x}|\hat{\boldsymbol{y}}\right)$. Since the BSC is discrete, $p\left(\hat{\boldsymbol{y}}|\boldsymbol{x};\boldsymbol{\varphi},\boldsymbol{\varepsilon}\right)$ is non-differentiable for $\boldsymbol{\varphi}$. Therefore, we sample $K$ noisy codewords $\hat{\boldsymbol{y}}$ for each input image $\boldsymbol{x}$ and use variational inference for Monte Carlo objectives (VIMCO) [37] to estimate $I_{\mathrm{VL}}\left(\boldsymbol{x};\hat{\boldsymbol{y}}\right)$ with low-variance gradients. The estimation $\hat{I}_{\mathrm{VL}}\left(\boldsymbol{x},\hat{\boldsymbol{y}};\boldsymbol{\varphi},\boldsymbol{\theta},\boldsymbol{\varepsilon}\right)$ can be expressed as [18], [19]

$$\hat{I}_{\mathrm{VL}}\left(\boldsymbol{x},\hat{\boldsymbol{y}};\boldsymbol{\varphi},\boldsymbol{\theta},\boldsymbol{\varepsilon}\right) =$$
$$\mathbb{E}_{p(\boldsymbol{x})} \mathbb{E}_{p\left(\hat{\boldsymbol{y}}^{(1):(K)}|\boldsymbol{x};\boldsymbol{\varphi},\boldsymbol{\varepsilon}\right)} \left[ \log \frac{1}{K} \sum_{i=1}^{K} p_{\boldsymbol{\theta}}\left(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}^{(i)}\right) \right], \quad (8)$$

where $\hat{\boldsymbol{y}}^{(i)}$ represents the $i$-th sample among $K$ samples.

To calculate (8), we need to know $p\left(\hat{\boldsymbol{y}}|\boldsymbol{x};\boldsymbol{\varphi},\boldsymbol{\varepsilon}\right)$ and $p_{\boldsymbol{\theta}}\left(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}\right)$ first. According to (1), (2), and (4), a Markov chain, $\boldsymbol{x} \to \boldsymbol{y} \to \hat{\boldsymbol{y}} \to \hat{\boldsymbol{x}}$ exists in the JSCC system. Thus, the joint probability $p\left(\boldsymbol{x},\hat{\boldsymbol{x}},\boldsymbol{y},\hat{\boldsymbol{y}}\right)$ can be modelled as

$$p\left(\boldsymbol{x},\hat{\boldsymbol{x}},\boldsymbol{y},\hat{\boldsymbol{y}}\right) = p\left(\boldsymbol{x}\right) p_{\boldsymbol{\varphi}}\left(\boldsymbol{y}|\boldsymbol{x}\right) p_{\boldsymbol{\varepsilon}}\left(\hat{\boldsymbol{y}}|\boldsymbol{y}\right) p_{\boldsymbol{\theta}}\left(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}\right), \quad (9)$$

where $p_{\boldsymbol{\varphi}}\left(\boldsymbol{y}|\boldsymbol{x}\right)$ is the conditional probability of $\boldsymbol{y}$ given $\boldsymbol{x}$, which is parameterized by the encoder neural network, and $p_{\boldsymbol{\varepsilon}}\left(\hat{\boldsymbol{y}}|\boldsymbol{y}\right)$ is the channel transition probability for BSC. Here, we consider the following two types of BSC:

- Single-channel: When the system utilizes single-carrier modulation such as binary phase shift keying (BPSK), the error probability of different bands is the same. This scenario is referred to as a single-channel scenario.
- Parallel-channel: When the system utilizes multicarrier modulation, e.g., orthogonal frequency division multiplexing (OFDM) to resist channel fading, the total available bandwidth is divided into non-overlapping bands, and the transmitted data stream will be divided into substreams and sent via parallel bands. In this case, each of the parallel bands has a different error probability. This scenario is referred to as a parallel-channel scenario. For instance, a parallel-channel scenario with 4 subchannels is shown in Fig. 3.

We assume that the system has $P$ parallel subchannels with equal bandwidth $\frac{M}{P}$, where $P \in \mathbb{N}^+$ and $M$ is the length of $\boldsymbol{y}$. Note that $P = 1$ represents the single-channel case. Denote the error probabilities of different bands as $\boldsymbol{\varepsilon} = \{\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_P\}$, the channel transition probability is

$$p_{\boldsymbol{\varepsilon}}(\hat{\boldsymbol{y}}|\boldsymbol{y}) = \prod_{k=1}^{P} \prod_{m=\frac{M}{P}(k-1)+1}^{\frac{M}{P}k} \varepsilon_k^{y_m \oplus \hat{y}_m}(1-\varepsilon_k)^{y_m \oplus \hat{y}_m \oplus 1},$$
(10)

where $\hat{y}_m$ represents the $m$-th element of $\hat{\boldsymbol{y}}$. Note that the proposed AIB-JSCC is also applicable to JSCC systems with arbitrary discrete memoryless channels (DMCs).

Let $f_{\boldsymbol{\varphi}}(\boldsymbol{x})$ denote the output of the encoder neural network when the input is $\boldsymbol{x}$. Since $\boldsymbol{y}$ is a binary codeword, we use the Bernoulli distribution to parameterize $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$. To reduce the redundancy between any two elements of $\boldsymbol{y}$, we assume that the elements in $\boldsymbol{y}$ are independent of each other, and $f_{\boldsymbol{\varphi}}(\boldsymbol{x})$ is treated as the parameters of this Bernoulli distribution, i.e., $\boldsymbol{y} = \mathrm{E}_{\boldsymbol{\varphi}}(\boldsymbol{x}) \sim \mathrm{Bern}(f_{\boldsymbol{\varphi}}(\boldsymbol{x}))$. Then, $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$ is

$$\begin{aligned} p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x}) &= \prod_{m=1}^{M} p_{\boldsymbol{\varphi}}(y_m|\boldsymbol{x}) \\ &= \prod_{m=1}^{M} (f_{\boldsymbol{\varphi}}(\boldsymbol{x}))^{y_m}(1 - f_{\boldsymbol{\varphi}}(\boldsymbol{x}))^{1-y_m}, \end{aligned}$$
(11)

where $y_m$ represents the $m$-th element of $\boldsymbol{y}$. The channel state information (CSI) is assumed to be perfectly estimated. Hence, both the encoder and the decoder know the accurate $\varepsilon$. We can compute $p(\hat{\boldsymbol{y}}|\boldsymbol{x}; \boldsymbol{\varphi}, \boldsymbol{\varepsilon})$ by marginalizing over $\boldsymbol{y}$ as

$$p(\hat{\boldsymbol{y}}|\boldsymbol{x}; \boldsymbol{\varphi}, \boldsymbol{\varepsilon}) = \sum_{\boldsymbol{y} \in \{0,1\}^M} p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x}) p_{\boldsymbol{\varepsilon}}(\hat{\boldsymbol{y}}|\boldsymbol{y}). \quad (12)$$

Then, $p(\hat{\boldsymbol{y}}|\boldsymbol{x}; \boldsymbol{\varphi}, \boldsymbol{\varepsilon})$ is formulated as:

$$p(\hat{\boldsymbol{y}}|\boldsymbol{x}; \boldsymbol{\varphi}, \boldsymbol{\varepsilon}) = \prod_{k=1}^{P} \prod_{m=\frac{M}{P}(k-1)+1}^{\frac{M}{P}k} (\xi_k(\boldsymbol{x}))^{\hat{y}_m}(1 - \xi_k(\boldsymbol{x}))^{1-\hat{y}_m},$$
(13)

where $\xi_k(\boldsymbol{x}) = f_{\boldsymbol{\varphi}}(\boldsymbol{x}) - 2f_{\boldsymbol{\varphi}}(\boldsymbol{x})\varepsilon_k + \varepsilon_k$. From (13), we can observe that $p(\hat{\boldsymbol{y}}|\boldsymbol{x}; \boldsymbol{\varphi}, \boldsymbol{\varepsilon})$ follows multivariate independent Bernoulli distribution with parameters $\xi_k(\boldsymbol{x})$.

Since $\boldsymbol{x}$ can be normalized to a real-value vector where each element value is within 0 and 1, we use a Gaussian distribution to model $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$ such that $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$ is differential with respect to $\boldsymbol{\theta}$. Let $g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})$ represent the output of the decoder neural network when the input of the decoder is $\hat{\boldsymbol{y}}$. We assume that the average of $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$ is $g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})$ [18], [19], i.e., $\hat{\boldsymbol{x}} = \mathrm{D}_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}}) \sim \mathcal{N}(g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}}), \boldsymbol{I})$, where $\mathcal{N}$ represents the Gaussian distribution. Then, we have

$$p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}}) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{\hat{x}_i - g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})_i}{2}\right), \quad (14)$$

where $\hat{x}_i$ is the $i$-th pixel of $\hat{\boldsymbol{x}}$, and $g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})_i$ is the corresponding pixel in $g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})$. Then $\hat{I}_{\mathrm{VL}}(\boldsymbol{x}, \hat{\boldsymbol{y}}; \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon})$ can be calculated by introducing (13) and (14) into (8).

## C. Upper Bound on $I(\boldsymbol{x}; \boldsymbol{y})$

Next, we derive the applicable form of $I(\boldsymbol{x}; \boldsymbol{y})$ in the considered system. Since $I(\boldsymbol{x}; \boldsymbol{y})$ is mathematically intractable, we minimize its upper bound instead. However, since we do not constrain the distribution of $\boldsymbol{y}$, the popular variational upper bound (VUB) [26], $KL(p(\boldsymbol{y}|\boldsymbol{x})|r(\boldsymbol{y}))$, cannot be used, where $r(\boldsymbol{y})$ is an approximation of $p(\boldsymbol{y})$. Therefore, we exploit another upper bound on mutual information called CLUB as [38]

$$\begin{aligned} I_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi}) &= \mathbb{E}_{p(\boldsymbol{x},\boldsymbol{y})}[\log p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})] \\ &\quad - \mathbb{E}_{p(\boldsymbol{x})}\mathbb{E}_{p(\boldsymbol{y})}[\log p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})]. \end{aligned}$$
(15)

Let $B$ denote the number of independent sample pairs $\{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^{B}$, where $\boldsymbol{x}^{(i)}$ represents the $i$-th image, and $\boldsymbol{y}^{(i)}$ represents the corresponding $i$-th codeword. Then $I_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$ can be estimated by the Monte Carlo method as:

$$\begin{aligned} \hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi}) &= \frac{1}{B}\sum_{i=1}^{B} \log p_{\boldsymbol{\varphi}}\left(\boldsymbol{y}^{(j)}|\boldsymbol{x}^{(i)}\right) \\ &\quad - \frac{1}{B^2}\sum_{i=1}^{B}\sum_{j=1}^{B} \log p_{\boldsymbol{\varphi}}\left(\boldsymbol{y}^{(j)}|\boldsymbol{x}^{(i)}\right). \end{aligned}$$
(16)

Since $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$ is the probability of a Bernoulli distribution, $\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$ is tractable and differentiable. Thus, instead of minimizing the true value of $I(\boldsymbol{x}; \hat{\boldsymbol{y}})$, we minimize $\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$.

Overall, by replacing $I(\boldsymbol{x}; \hat{\boldsymbol{y}})$ and $I(\boldsymbol{x}; \boldsymbol{y})$ in (6) with $\hat{I}_{\mathrm{VL}}(\boldsymbol{x}, \hat{\boldsymbol{y}}; \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon})$ in (8) and $\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$ in (16), respectively, we can obtain a tractable and differential form of IB objective for the JSCC system as:

$$\max_{\boldsymbol{\varphi}, \boldsymbol{\theta}} \left[\hat{I}_{\mathrm{VL}}(x, \hat{\boldsymbol{y}}; \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}) - \beta\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})\right]. \quad (17)$$

Even though (17) can be used for training, the value of $\beta$ needs to be carefully optimized, which controls the trade-off between the compression level and the reconstruction quality. Therefore, in Section IV, we further propose an adaptive IB algorithm to determine the appropriate value of $\beta$.

## IV. ADAPTIVE IB ALGORITHM

This section first proposes an adaptive IB algorithm to select appropriate value of $\beta$ according to the distortion of reconstruction during the training process. We then describe the whole training process of AIB-JSCC which combines the proposed IB objective and the adaptive IB algorithm.

### A. Adaptive IB Algorithm

Since the values of $I(\boldsymbol{x}; \hat{\boldsymbol{y}})$ and $I(\boldsymbol{x}; \boldsymbol{y})$ change during the training process, it is necessary to alter $\beta$ accordingly so as to balance $I(\boldsymbol{x}; \hat{\boldsymbol{y}})$ and $I(\boldsymbol{x}; \boldsymbol{y})$. To adaptively determine the value of $\beta$ in each epoch, we propose a proportional-integral-differential (PID) control based algorithm, which determines the current value of $\beta$ by analyzing the past errors and

predicting future errors. The discrete form of PID controller can be expressed as [39]

$$\beta[w] = K_p e[w] - K_i \sum_{k=0}^{w-1} e[k] - K_d \left( e[w] - e[w-1] \right),$$
(18)

where $\beta[w]$ is the output of the controller at time $w$. $K_p$, $K_i$, $K_d$ and error $e[w]$ are the proportional gain, the integral gain, the differential gain and the difference between the actual value and the desired value at time $w$, respectively. In addition, $K_p e[w]$ is the proportional (P) term, which responds to the change of error quickly and provides a global control proportional to the error; $K_i \sum_{k=0}^{w-1}$ is the integral (I) term, which continues to increase as long as the error is greater than 0 and is used to eliminate steady-state errors; $K_d \left( e[w] - e[w-1] \right)$ is the differential (D) term, which can reduce the overshoot and improve the system's stability and transient response [39]. The PID controller continuously calculates error $e[w]$ and the weighted sum of these three terms, and then applies a correction on the system to reduce the error $e[w]$. We employ (18) to adjust $\beta$ at the end of each epoch.

However, before applying (18) to the AIB-JSCC system, an upper bound of $\beta$ must to derived. This is because if $\beta$ is excessively large, $I(\boldsymbol{x};\boldsymbol{y})$ will dominate the loss function, and $\boldsymbol{y}$ will aggressively discard information of $\boldsymbol{x}$, leading to the loss of useful information. Considering an extreme case when $\beta$ approaches positive infinity, $I(\boldsymbol{x};\boldsymbol{y})$ will approach 0, and in this case, the global optimal encoder distribution may be $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x}) = p(\boldsymbol{y})$. That means $\boldsymbol{y}$ becomes independent of $\boldsymbol{x}$. In this case, $\boldsymbol{y}$ and $\hat{\boldsymbol{y}}$ contain no information about $\boldsymbol{x}$, i.e. $I(\boldsymbol{x};\hat{\boldsymbol{y}}) = I(\boldsymbol{x};\boldsymbol{y}) = 0$, and reconstructing $\boldsymbol{x}$ from $\hat{\boldsymbol{y}}$ becomes infeasible. Therefore, it is necessary to limit $\beta$ below an upper bound before applying PID controller, as shown in the following lemma.

**Lemma 1.** The condition that $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x}) = p(\boldsymbol{y})$ is not a local optimum for the IB objective is [40]

$$\beta < \beta_{\max} = \sup_{\boldsymbol{x} \to \boldsymbol{y} \to \hat{\boldsymbol{y}}} \frac{I(\boldsymbol{x};\hat{\boldsymbol{y}})}{I(\boldsymbol{x};\boldsymbol{y})}.$$
(19)

*Proof.* See Appendix. □

According to Lemma 1, we derive the estimated upper bound on $\beta$ of the proposed AIB-JSCC in the following theorem.

**Theorem 1.** For $B$ pairs $\left\{ \left( \boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}, \hat{\boldsymbol{y}}^{(i)}, \hat{\boldsymbol{x}}^{(i)} \right) \right\}_{i=1}^{B}$, the estimated upper bound on $\beta$ is

$$\beta_{\max} = \frac{\hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}}}{\hat{I}_{\boldsymbol{x},\boldsymbol{y}}},$$
(20)

where

$$\hat{I}_{\boldsymbol{x},\boldsymbol{y}} = \sum_{m=1}^{M} H \left( \frac{1}{B} \sum_{i=1}^{B} p \left( y_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right)$$
$$- \frac{1}{B} \sum_{m=1}^{M} \sum_{i=1}^{B} H \left( p \left( y_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right),$$
(21)

and

$$\hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}} = \sum_{m=1}^{M} H \left( \frac{1}{B} \sum_{i=1}^{B} p \left( \hat{y}_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right)$$
$$- \frac{1}{B} \sum_{m=1}^{M} \sum_{i=1}^{B} H \left( p \left( \hat{y}_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right).$$
(22)

*Proof.* At the end of each epoch, $I(\boldsymbol{x};\hat{\boldsymbol{y}})$ and $I(\boldsymbol{x};\boldsymbol{y})$ are fixed since neural networks of the encoder and decoder are fixed. Therefore, according to Lemma 1, $\beta_{\max} = \frac{I(\boldsymbol{x};\hat{\boldsymbol{y}})}{I(\boldsymbol{x};\boldsymbol{y})}$. We then estimate $I(\boldsymbol{x};\hat{\boldsymbol{y}})$ and $I(\boldsymbol{x};\boldsymbol{y})$ according to the definition of the mutual information, i.e., $I(\boldsymbol{x};\hat{\boldsymbol{y}}) = H(\hat{\boldsymbol{y}}) - H(\hat{\boldsymbol{y}}|\boldsymbol{x})$, and $I(\boldsymbol{x};\boldsymbol{y}) = H(\boldsymbol{y}) - H(\boldsymbol{y}|\boldsymbol{x})$. We estimate $H(\boldsymbol{y})$ and $H(\boldsymbol{y}|\boldsymbol{x})$ separately. To obtain $H(\boldsymbol{y})$, we calculate the probability of the $m$-th element of $\boldsymbol{y}$, $p(y_m)$, by,

$$p(y_m) = \frac{1}{B} \sum_{i=1}^{B} p \left( y_m^{(i)} | \boldsymbol{x}^{(i)} \right),$$
(23)

where $y_m^{(i)}$ represents the $m$-th element of the codeword of the $i$-th input image $\boldsymbol{x}^{(i)}$. Since the elements in $\boldsymbol{y}$ are assumed to be independent, the entropy of $\boldsymbol{y}$ is equal to the sum of the entropies of all elements. Besides, we assume $p(\boldsymbol{x}) = \frac{1}{B}$, and we have

$$H(\boldsymbol{y}) = \sum_{m=1}^{M} H(y_m) \approx \hat{H}(\boldsymbol{y})$$
$$= \sum_{m=1}^{M} H \left( \frac{1}{B} \sum_{i=1}^{B} p \left( y_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right).$$
(24)

Substituting (11) into (24), $H(\boldsymbol{y})$ can be calculated. To calculate $H(\boldsymbol{y}|\boldsymbol{x})$, we use the assumption $p(\boldsymbol{x}) = \frac{1}{B}$ again, and we have

$$H(\boldsymbol{y}|\boldsymbol{x}) \approx \hat{H}(\boldsymbol{y}|\boldsymbol{x})$$
$$= \frac{1}{B} \sum_{m=1}^{M} \sum_{i=1}^{B} H \left( p \left( y_m^{(i)} | \boldsymbol{x}^{(i)} \right) \right).$$
(25)

Therefore, by exploiting $I(\boldsymbol{x};\boldsymbol{y}) = H(\boldsymbol{y}) - H(\boldsymbol{y}|\boldsymbol{x})$, we have

$$I(\boldsymbol{x};\boldsymbol{y}) \approx \hat{I}_{\boldsymbol{x},\boldsymbol{y}} = \hat{H}(\boldsymbol{y}) - \hat{H}(\boldsymbol{y}|\boldsymbol{x}).$$
(26)

Similar to (24), (25). and (26), we can estimate $I(\boldsymbol{x};\hat{\boldsymbol{y}})$ as

$$I(\boldsymbol{x};\hat{\boldsymbol{y}}) \approx \hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}} = \hat{H}(\hat{\boldsymbol{y}}) - \hat{H}(\hat{\boldsymbol{y}}|\boldsymbol{x}).$$
(27)

Given (26) and (27) we can obtain the upper bound as (20). This completes the proof. □

From Theorem 1, we can observe that both $p(\boldsymbol{y}|\boldsymbol{x})$ and $p(\hat{\boldsymbol{y}}|\boldsymbol{x})$ affect $\beta_{\max}$, which can be calculated via (11) and (13) in our designed system.

To ensure the relevance between the source $\boldsymbol{x}$ and the codeword $\boldsymbol{y}$, we further add the approximated upper bound on $\beta$, $\beta_{\max}[w]$ on the basis of PID algorithm to constrain the range of $\beta$. Note that the optimal value of $\beta$ to minimize MSE may not be 0 since we jointly optimize the transmission rate and the distortion of JSCC. In the ideal case, the transmission distortion can be reduced close to 0 under a certain compression ratio. Therefore, we treat the MSE between the original image in the validation set and the corresponding reconstructed

**Algorithm 1** Adaptive IB Algorithm

**Input:** Encoder $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$; Decoder $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$; MSE at $(w-1)$-th epoch and $w$-th epoch: MSE$[w]$ and MSE$[w-1]$, Coefficients $K_p$, $K_i$ and $K_d$; Minimal value of hyperparameter $\beta$ : $\beta_{\min}$.

**Output:** The hyperparameter used in $w$-th epoch: $\beta[w]$.

**Initialization:** $I[0] = 0$; MSE $[0] = 0$.

1: $P[w] \leftarrow K_p \text{MSE}[w]$;
2: $I[w] \leftarrow I[w-1] + K_i \text{MSE}[w]$;
3: $D[w] \leftarrow K_d (\text{MSE}[w] - \text{MSE}[w-1])$;
4: Calculate $\beta_{\max}[w]$ according to (20);
5: $\beta_{\text{adp}}[w] \leftarrow \beta_{\max}[w] + P[w] - I[w] - D[w]$;
6: $\beta[w] \leftarrow \text{clamp}(\beta_{\text{adp}}[w], \beta_{\min}, \beta_{\max}[w])$;

---

image at the $w$-th epoch as $e[w]$. Then, by applying the PID algorithm, $\beta$ will change in the direction of reducing MSE. The proposed formula of adaptive $\beta$ at the $w$-th epoch is

$$\beta_{\text{adp}}[w] = \beta_{\max}[w] + K_p \text{MSE}[w] - K_i \sum_{k=0}^{w-1} \text{MSE}[k] \qquad (28)$$
$$- K_d (\text{MSE}[w] - \text{MSE}[w-1]),$$

where $\beta_{\max}[w] = \frac{\hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}}[w]}{\hat{I}_{\boldsymbol{x},\boldsymbol{y}}[w]}$, $\hat{I}_{\boldsymbol{x},\boldsymbol{y}}[w]$ and $\hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}}[w]$ are $\hat{I}_{\boldsymbol{x},\boldsymbol{y}}$ and $\hat{I}_{\boldsymbol{x},\hat{\boldsymbol{y}}}$ at the $w$-th epoch, respectively, MSE $[w]$ is the average MSE at the $w$-th epoch and it is expressed as

$$\text{MSE}[w] = \frac{1}{V} \frac{1}{N} \sum_{i=1}^{V} \sum_{j=1}^{N} \left( x_j^{(i)}[w] - \hat{x}_j^{(i)}[w] \right)^2, \qquad (29)$$

with $V$ being the number of the images in the validation set, $x_j^{(i)}[k]$ being the $j$-th pixel in the $i$-th transmitted image $\boldsymbol{x}^{(i)}$ recovered at the $k$-th epoch, and $\hat{x}_j^{(i)}[k]$ being the corresponding pixel in the corresponding reconstructed image $\hat{\boldsymbol{x}}^{(i)}$ at the $w$-th epoch.

After training, $I(\boldsymbol{x};\hat{\boldsymbol{y}})$ and $I(\boldsymbol{x};\boldsymbol{y})$ slightly fluctuates in a small range, and the balance between them is nearly fixed. $\beta$ should converge to a certain minimal value. In general, the minimal value of $\beta$ is larger than 0 since $\beta = 0$ means ignoring the compression term $I(\boldsymbol{x};\boldsymbol{y})$. Therefore, we constrain $\beta$ larger than a minimum value, $\beta_{\min}(>0)$. Then, $\beta$ at $w$-th epoch can be expressed as:

$$\beta[w] = \text{clamp}(\beta_{\text{adp}}[w], \beta_{\min}, \beta_{\max}[w]), \qquad (30)$$

where $\text{clamp}(x, \min, \max)$ represents clamping $\boldsymbol{x}$ between $\min$ and $\max$ ($\min \leq \max$).

From (6), the importance of $I(\boldsymbol{x};\boldsymbol{y})$ in the loss function decreases as $\beta$ increases. At the beginning, $\boldsymbol{y}$ contains abundant redundancy due to imperfect map from the source information to the transmitted codewords. Therefore, in the initial stages, we must set a relatively large $\beta$ to squeeze more redundant information in $\boldsymbol{y}$. As the training processes, the value of $\beta$ should decrease since the corresponding redundancy information in $\boldsymbol{y}$ gradually decreases. The value of $\beta$ will finally converge to a constant when the proposed AIB-JSCC achieves the optimal trade-off between the reconstruction quality and the compression ratio. Therefore, we adjust the coefficients

**Algorithm 2** AIB-JSCC

**Input:** Dataset($\mathcal{X}$) to be compressed; Channel error probability $\varepsilon$; Hyperparameter $\beta$.

**Output:** Learned encoder $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$ and decoder $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$.

1: Initialize the parameters of encoder $p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$, the parameters of decoder $p_{\boldsymbol{\theta}}(\hat{\boldsymbol{x}}|\hat{\boldsymbol{y}})$; $i = 1$.
2: **while** not converge **do**
3:     Sample $B$ samples from Dataset: $\boldsymbol{x} \sim p(\boldsymbol{x})$;
4:     Sample a codeword $\boldsymbol{y} \sim p_{\boldsymbol{\varphi}}(\boldsymbol{y}|\boldsymbol{x})$ for each $\boldsymbol{x}$;
5:     Sample $K$ noisy codewords $\hat{\boldsymbol{y}} \sim p(\hat{\boldsymbol{y}}|\boldsymbol{x};\boldsymbol{\varphi},\varepsilon)$ for each $\boldsymbol{x}$;
6:     Calculate $\hat{I}_{\text{VL}}(\boldsymbol{x},\hat{\boldsymbol{y}};\boldsymbol{\varphi},\boldsymbol{\theta},\varepsilon)$ according to (8);
7:     Calculate $\hat{I}_{\text{CLUB}}(\boldsymbol{x},\boldsymbol{y};\boldsymbol{\varphi})$ according to (16);
8:     Update $\boldsymbol{\varphi}$ and $\boldsymbol{\theta}$ according to (17);
9:     **if** an epoch of training finishes **then**
10:         Calculate MSE [i] according to (29);
11:         Update $\beta[i]$ according to Algorithm 1;
12:     **end if**
13:     $i \leftarrow i + 1$;
14: **end while**

---

$K_p$, $K_i$ and $K_d$ to let $\beta$ gradually decrease from its upper bound as the training processes. The adaptive IB algorithm is summarized in Algorithm 1.

*B. Training Process of AIB-JSCC*

The architecture of the AIB-JSCC system is shown in Fig. 2. The encoder first extracts the information of the input image as a feature map according to a feature extractor which consists of convolutional neural networks (CNN) or fully connected (FC) layer. Then, to control the length of $\boldsymbol{y}$, an FC layer is used to turn the feature map into a $M$-dimensional vector $f_{\boldsymbol{\varphi}}(\boldsymbol{x})$. The codeword $\boldsymbol{y}$ is sampled according to $y \sim \text{Bern}(f_{\boldsymbol{\varphi}}(\boldsymbol{x}))$. At the receiver, the noisy codeword $\hat{\boldsymbol{y}}$ is first passed into an FC layer and then reshaped into a feature map. The feature map is upsampled to the same dimension as $\boldsymbol{x}$ to obtain $g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}})$. The recovered image $\hat{\boldsymbol{x}}$ is generated according to $\hat{\boldsymbol{x}} \sim N(g_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}}), \boldsymbol{I})$. At the output layer of the decoder, an activation function is used to transform the pixel values in $\hat{\boldsymbol{x}}$ to $[0, 1]$. We then multiply $\hat{\boldsymbol{x}}$ by 255 and round the resulting values to ensure that the pixel values are discrete and fall between $[0, 255]$. (17) is used as the loss function to train the encoder and the decoder for image transmission jointly. We use the mini-batch gradient descent method [41] to optimize the parameters. To guarantee that each image in the batch has an equal probability of being selected for updating the parameters, we have $p(\boldsymbol{x}) = \frac{1}{B}$. Then, the distribution of $\boldsymbol{y}$ can be obtained shown in (11), and the distribution of $\hat{\boldsymbol{y}}$ can be obtained shown in (13). We sample one codeword $\boldsymbol{y}$ for each $\boldsymbol{x}$, and have $B$ pairs $\left\{ \left( \boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)} \right) \right\}_{i=1}^{B}$. Then, the transmission rate, $\hat{I}_{\text{CLUB}}(\boldsymbol{x},\boldsymbol{y};\boldsymbol{\varphi})$ is calculated based on (16). To calculate $\hat{I}_{\text{VL}}(\boldsymbol{x},\hat{\boldsymbol{y}};\boldsymbol{\varphi},\boldsymbol{\theta},\varepsilon)$, we further sample $K$ noisy codewords $\hat{\boldsymbol{y}}$ for each $\boldsymbol{x}$, and the total number of $\hat{\boldsymbol{y}}$ is $B \times K$. According to (8), the distortion term $\hat{I}_{\text{VL}}(\boldsymbol{x},\hat{\boldsymbol{y}};\boldsymbol{\varphi},\boldsymbol{\theta},\varepsilon)$ can be calculated. Finally, $\boldsymbol{\varphi}$ and $\boldsymbol{\theta}$ are updated according to (17). At the end of each epoch, Algorithm 1 is applied to update

TABLE II
SYSTEM PARAMETERS.

| Parameters | Value | | | |
|---|---|---|---|---|
| Datasets | MNIST | CIFAR10 | SVHN | Omniglot |
| Codewords length $M$ | 100 | 400,450, 500,550,600 | 500 | 200 |
| Channel error probability $\varepsilon$ | Single-channel: 0.1, 0.2, 0.3, 0.4 | | | |
| | Parallel-channel: shown as (32) | | | |
| Subchannel number $P$ | $2, 4, 5$ | | | |
| $\beta$ ( $\beta_{\min}$ ) | 0.01 | | | 0.001 |
| $K_p$ | 0.001 | | | |
| $K_i$ | $-0.001$ | $-0.0001$ | | |
| $K_d$ | $-0.001$ | | | |
| Batchsize $B$ | 300 | | | |
| Training Epoch | 500 | | | |
| Learning Rate | 0.001 | | | |
| Regularization coefficient | 0.0001 | | | |

TABLE III
PSNR OF NECST VS. IABF VS. IB-JSCC VS. AIB-JSCC.

| Datasets | Methods | PSNR under different error probabilities | | | |
|---|---|---|---|---|---|
| | | 0.1 | 0.2 | 0.3 | 0.4 |
| MNIST | NECST | 17.348 | 15.411 | 13.581 | 12.104 |
| | IABF | 17.721 | 15.513 | 13.735 | 12.264 |
| | IB-JSCC | 17.801 | 15.724 | 13.741 | 12.408 |
| | AIB-JSCC | **17.837** | **15.725** | **13.751** | **12.411** |
| Omniglot | NECST | 15.017 | 13.955 | 12.959 | 12.1409 |
| | IABF | 15.117 | 13.928 | 13.039 | 12.166 |
| | IB-JSCC | 15.158 | 14.015 | 13.04 | 12.203 |
| | AIB-JSCC | **15.161** | **14.03** | **13.052** | **12.213** |
| CIFAR10 | NECST | 16.864 | 16.158 | 15.35 | 14.163 |
| | IABF | 17.442 | 16.391 | 15.673 | 14.219 |
| | IB-JSCC | 17.455 | 16.68 | 15.792 | 14.247 |
| | AIB-JSCC | **17.513** | **16.748** | **15.809** | **14.282** |

$\beta$. The coefficients $K_p$, $K_i$, and $K_d$ are adjusted to obtain proper $\beta$. The updated $\beta$ is then used in the loss function of the next epoch. The model with the lowest average MSE on valid dataset during training is stored. After training, the AIB-JSCC system can reduce the transmitted data by minimizing $\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$ while improving the reconstruction quality by maximizing $\hat{I}_{\mathrm{VL}}(\boldsymbol{x}, \hat{\boldsymbol{y}}; \boldsymbol{\varphi}, \boldsymbol{\theta}, \varepsilon)$. The whole training procedure of the proposed AIB-JSCC is summarized in Algorithm 2.

## V. EXPERIMENTAL RESULTS

In this section, we provide extensive experiments to validate our designed system. The experiments are carried on the following datasets: MNIST [42], Omniglot [43], CIFAR10 [44] and street view housing numbers tsiscon (SVHN) [45] to account for different image sizes and colors. For comparison purposes, we choose IABF and NECST in [19] and [18]



Fig. 4. Visual comparison between IABF and IB-JSCC.

and classical SSCC schemes as baselines. Specifically, for SSCC schemes, we employ three industry-standard source encoders: JPEG [46], JPEG2000 [47] and WebP [48], and BPG [49], combined with an ideal capacity-achieving channel code (marked as "JPEG + Capacity", "JPEG2000 + Capacity" and "WebP + Capacity", and "BPG + Capacity", respectively). We do not compare with LDPC coding as we are often unable to obtain valid image files after LDPC decoding. To make a fair comparison, the settings and structure of the encoder and decoder neural networks used in AIB-JSCC are the same as those used in IABF. The system parameters are shown in Table II. In line with the baseline and references [18], [19], we choose similar parameters of the neural network and the optimizer. For Monte Carlo estimation of $I_{\mathrm{VL}}(\boldsymbol{x}; \hat{\boldsymbol{y}})$, we utilize 5 samples [18] [19]. The Monte Carlo estimates of $\hat{I}_{\mathrm{CLUB}}(\boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\varphi})$, $\mathrm{H}(\boldsymbol{y})$, $\mathrm{H}(\boldsymbol{y}|\boldsymbol{x})$ and $\mathrm{H}(\hat{\boldsymbol{y}})$ use 300 samples per batch. We choose the best $K_p$ and $K_d$ from the set $\left\{10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\right\}$ that provide the best performance. We use widely-used image quality metrics, MSE and peak signal-to-noise ratio (PSNR), to measure the performance of AIB-JSCC and IABF [19]. In the single-channel scenario, we compare the reconstruction and compression ability of AIB-JSCC with baselines, and the robustness and complexity of AIB-JSCC are also discussed. In the parallel-channel scenario, we present the results of the reconstruction error, the distribution of neuron weights and the visual reconstructions to illustrate that AIB-JSCC can adaptively allocate elements to parallel channels according to their channel state information. The above experiments are implemented for the BSC. We also compare the reconstruction error of AIB-JSCC and IABF when the channel is the high-order DMC to demonstrate the effectiveness of AIB-JSCC.

### A. Single-channel Scenario

*1) Reconstruction capabilitiy:* Table III shows the PSNR of different schemes on various datasets under different error probabilities, where IB-JSCC stands for the degenerate AIB-JSCC with fixed $\beta$. The definition of PSNR is

$$\mathrm{PSNR} = 10\log_{10}\left(\frac{(2^n - 1)^2}{\mathrm{mse}(\boldsymbol{x}, \hat{\boldsymbol{x}})}\right), \qquad (31)$$

where $n$ is the number of bits that each image pixel uses, $\mathrm{mse}(\boldsymbol{x}, \hat{\boldsymbol{x}})$ is MSE between $\boldsymbol{x}$ and $\hat{\boldsymbol{x}}$. In particular, we fix the length of $\boldsymbol{y}$ and calculate the average MSE and PSNR over the test sets. From Table III, we can observe that AIB-JSCC is always superior to IB-JSCC and IABF in terms of MSE

TABLE IV
CLASSIFICATION ACCURACY OF IMAGES RECOVERED BY IABF AND
AIB-JSCC

| Classifiers | Methods | Acc under different error probabilities | | | |
|---|---|---|---|---|---|
| | | **0.1** | **0.2** | **0.3** | **0.4** |
| MLP | IABF | 0.932 | 0.817 | 0.637 | 0.331 |
| | AIB-JSCC | **0.937** | **0.881** | **0.692** | **0.386** |
| SVM | IABF | 0.932 | 0.821 | 0.619 | 0.312 |
| | AIB-JSCC | **0.942** | **0.884** | **0.694** | **0.358** |
| DT | IABF | 0.51 | 0.391 | 0.297 | 0.177 |
| | AIB-JSCC | **0.564** | **0.469** | **0.347** | **0.2** |
| RF | IABF | 0.673 | 0.522 | 0.288 | 0.176 |
| | AIB-JSCC | **0.708** | **0.549** | **0.308** | **0.181** |



Fig. 6. MSE of IABF and AIB-JSCC with different $M$.



Fig. 5. The value of hyperparameter $\beta[w]$ with respect to training epoch.



Fig. 7. The additional number of bits need by SSCC.

and PSNR, which validates the effectiveness of the proposed IB objective and the adaptive IB algorithm. From Table III, we can also observe that AIB-JSCC and IB-JSCC can reduce MSE and increase PSNR more on the RGB dataset CIFAR10 than on the greyscale datasets MNIST and Omniglot. This is because AIB-JSCC can extract information more precisely with the guidance of the proposed IB objective, thus recovering complex images better.

Figure 4 shows the visual reconstructions of AIB-JSCC and IABF on the SVHN dataset where $\varepsilon = 0.1$ and $M = 500$. From Fig. 4, we can observe that images recovered by AIB-JSCC are closer to original ones than those recovered by IABF, and the numbers in images recovered by AIB-JSCC can be distinguished more easily than IABF. For example, the first image recovered by AIB-JSCC can be clearly recognized as 5 while the one recovered by IABF may be incorrectly recognized as 9 or 4. This implies that AIB-JSCC can preserve more semantic information than IABF. This is due to the fact that compared with IABF, AIB-JSCC can preserve useful information as well as discard useless information which may lead to semantic mistakes. In consequence, AIB-JSCC has better visual reconstruction quality.

Table IV shows the classification accuracy (Acc) of the images reconstructed by IABF and AIB-JSCC with respect to different error probabilities. In particular, 4 different classifiers, multilayer perceptron (MLP), support vector machines (SVM), decision trees (DT) and random forests (RF) are trained with the raw MNIST train set, and tested with the images reconstructed by IABF and AIB-JSCC from the MNIST test set where $\varepsilon = 0.1$. From Table IV, we can observe that the classification accuracy of the images recovered by AIB-JSCC is always higher than IABF over different error probabilities. This implies that AIB-JSCC can preserve more semantic information useful for downstream task. This is because the proposed IB objective preserves information as well as discards information. Hence, AIB-JSCC can extract information more precisely.

Figure 5 shows the trend of AIB-JSCC's hyperparameter $\beta[w]$ with respect to training epoch on CIFAR10 dataset under different error probabilities. From Fig. 5, we can observe that $\beta[w]$ gradually decreases to a minimal value $\beta_{\min}$ when training processes. This is due to the fact that when the training
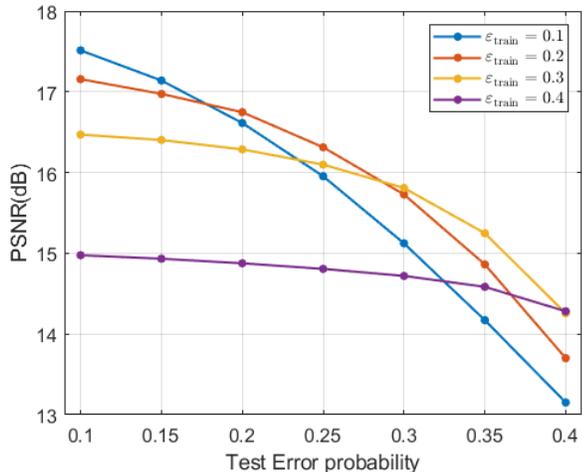
Fig. 8. The PSNR of AIB-JSCC with various train and test error probabilities.

TABLE VI
THE INFERENCE TIME OF AIB-JSCC AND SSCC.

| Methods | AIB-JSCC | BPG | WebP | JPEG2000 | JPEG |
|---|---|---|---|---|---|
| Inference time (ms) | 0.8 | 109 | 1.632 | 1.387 | 0.911 |

TABLE VII
MSE OF AIB-JSCC UNDER DIFFERENT PARALLEL-CHANNEL SCENARIOS

| Scenario | Average of $\varepsilon$ | Datasets | | |
|---|---|---|---|---|
| | | MNIST | Omniglot | CIFAR10 |
| $2-$ch | 0.051 | **9.616** | **23.538** | **45.069** |
| $4-$ch | 0.138 | 12.398 | 27.599 | 52.834 |
| $5-$ch | 0.213 | 14.075 | 28.927 | 54.021 |
| single-channel | 0.1 | 12.902 | 23.89 | 54.464 |
| | 0.2 | 20.784 | 31.092 | 64.969 |

TABLE V
THE NUMBER OF NETWORK PARAMETERS AND THE INFERENCE TIME OF
IABF AND AIB-JSCC.

| Datasets | Methods | parameters number ($\times 10^5$) | Inference time (ms) |
|---|---|---|---|
| MNIST | IABF | 11.86 | 0.4 |
| | AIB-JSCC | **11.345** | 0.4 |
| Omniglot | IABF | 13.36 | 0.743 |
| | AIB-JSCC | **12.345** | **0.372** |
| CIFAR10 | IABF | 5.679 | 1 |
| | AIB-JSCC | **3.164** | **0.8** |

processes, $I(x;\hat{y})$ increases and $I(x;y)$ decreases, and to keep balance between $I(x;\hat{y})$ and $I(x;y)$, the proposed adaptive IB algorithm decreases $\beta[w]$ to reduce the proportion of $I(x;\hat{y})$ in the loss function. From Fig. 5, we can also observe that the value of $\beta[w]$ reduces as the error probability increases. This is because when the channel error probability increases, we need to add more redundancy to the transmitted codeword $y$. In AIB-JSCC, this is achieved by increasing the distortion term $I(x;\hat{y})$. Since there is a Markov chain relationship $x \to y \to \hat{y} \to \hat{x}$, $I(x;\hat{y}) \leq I(x;y) \leq H(y)$, and so maximizing $I(x;\hat{y})$ can essentially increase $I(x;y)$ and $H(y)$, thus increasing the redundancy in $y$. This implies that the proposed adaptive IB algorithm is able to adjust $\beta[w]$ according to $I(x;\hat{y})$, $I(x;y)$ and the error probability.

*2) Compression Capability:* In this section, we denote the length of $y$, $M$, used by AIB-JSCC as the number of bits in order to compare the compression capability with other baselines. Note that $M$ bits is an upper bound on the transmission rate $I(x;y)$. Figure 6 shows the reconstruction MSE of IABF and AIB-JSCC with different $M$ on CIFAR10 under different error probabilities. In particular, to guarantee fairness, we use the results of IABF present in [19] in order to prevent the performance reduction caused by improper hyperparameter selection. From Fig. 6, we can observe that to obtain similar MSE, AIB-JSCC requires 15, 100, 70, 20 fewer bits than IABF when $\varepsilon = 0.1, 0.2, 0.3, 0.4$. This implies

that AIB-JSCC can reduce more than $20\%$ transmission rate compared with IABF. The $20\%$ gain stems from the fact that AIB-JSCC simultaneously minimizes the distortion and the transmission rate thus reducing the transmission rate to achieve a similar reconstruction error.

Figure 7 shows the additional number of bits that SSCC schemes need to achieve similar MSE on SVHN, compared with AIB-JSCC. From Fig. 7, we can observe that SSCC schemes need more bits than AIB-JSCC at all datasets and error probabilities. We can also observe that although BPG + Capacity needs fewer bits than the other three SSCC schemes, BPG + Capacity still needs more bits than AIB-JSCC for all datasets and error probabilities. When the error probability $\varepsilon$ increases, the additional required number of bits will increase. When $\varepsilon = 0.4$, AIB-JSCC only needs around $4\%$ JPEG needs. The $4\%$ gains stem from the fact that SSCC schemes are designed to be optimized for squared error with hand-selected constraints [21], [50], [51] while AIB-JSCC jointly trains the encoder and decoder by maximizing $I(x;\hat{y})$ and minimizing $I(x,y)$ thus preserving information precisely with lower transmission rate.

*3) Complexity and robustness:* The most computationally costly operations in the network are the convolutions/deconvolutions and the FC layers, as they involve multiplications and additions. The computational cost of a single convolutional layer is $H \times W \times K \times K \times C_i \times C_o$, where $K$ is the filter size, $C_o$ is the number of output channels, $C_i$ is the number of input channels and $H \times W$ is the size of the feature map. The computational cost of an FC layer is $(2I - 1)O$, where $I$ is the input vector dimension and $O$ is the output vector dimension. Only the width and height of the feature map and the vector dimension depend on the image dimensions, and all other factors are constant and independent of the image size. Thus, the computational complexity of the proposed scheme is $O(I_H \times I_W)$, where $I_H$ and $I_W$ are the width and height of the input image.

Table V shows the number of encoder and decoder network parameters and the inference time of IABF and AIB-JSCC. From Table V, we can observe that AIB-JSCC has fewer
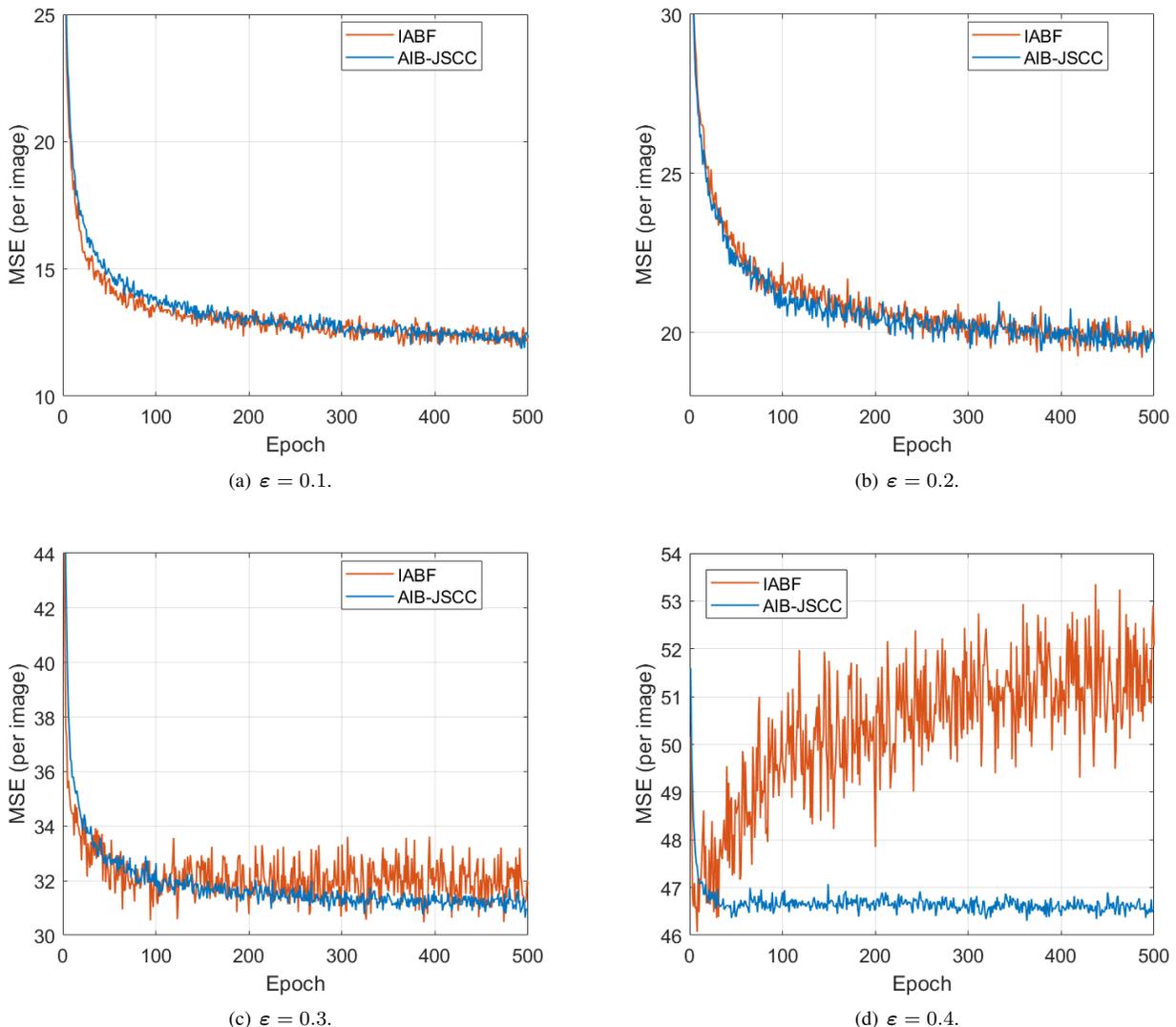
Fig. 9. MSE of IABF and AIB-JSCC on MNIST dataset. The error is calculated on validation set during training.

parameters and needs less inference time. Specifically, AIB-JSCC can reduce 45% parameters on CIFAR10 and 50% inference time on Omniglot. These gains stem from the simple network structure of AIB-JSCC, which makes the computational complexity of AIB-JSCC lower than that of IABF.

Table VI shows the inference time of AIB-JSCC, JPEG, JPEG2000, WebP and BPG on SVHN. From Table VI, we can observe that AIB-JSCC achieves lower inference time compared to all considered source coding schemes. Furthermore, it is worth noting that SSCC necessitates iterative channel decoding to attain optimal error correction capability [52]–[54]. As a result, the time required by SSCC is significantly higher than that of source coding. Consequently, in comparison to practical SSCC, AIB-JSCC is expected to yield superior time savings, which are not entirely reflected in Table VI.

Figure 8 shows the PSNR of AIB-JSCC on CIFAR10 when there is an estimation error on the channel error probability. From Fig. 8, we can observe that when $\varepsilon_{\text{test}}$ drops below $\varepsilon_{\text{train}}$, the performance does not saturate immediately. When

$\varepsilon_{\text{test}}$ increases beyond $\varepsilon_{\text{train}}$, AIB-JSCC exhibits a graceful degradation of the reconstruction quality. This is because AIB-JSCC uses the channel conditions in the loss function and enables the learned codewords to resist channel interference. Hence, the codewords extracted by AIB-JSCC is robust to different error probabilities.

Figure 9 shows the changes of validation reconstruction MSE with respect to training time steps for IABF and AIB-JSCC. From Fig. 9, we can observe that the trends of IABF and AIB-JSCC are similar when error probability $\varepsilon$ is 0.1 and 0.2. As the error probability gets larger, AIB-JSCC converges more stably than IABF. For example, in Fig. 9(d), when $\varepsilon$ is 0.4, there is severe overfitting in IABF while AIB-JSCC still converges stably. This is because AIB-JSCC avoids overfitting according to minimizing $I(\boldsymbol{x}; \boldsymbol{y})$ by neural network. Therefore, AIB-JSCC is more robust than IABF.

Figure 10 shows the 2-dimensional projections of the noisy codewords extracted from MNIST with different test error probabilities, and each color represents a number, i.e. a type

(a) $\varepsilon_{\text{test}} = 0.1$.

(b) $\varepsilon_{\text{test}} = 0.2$.

(c) $\varepsilon_{\text{test}} = 0.3$.
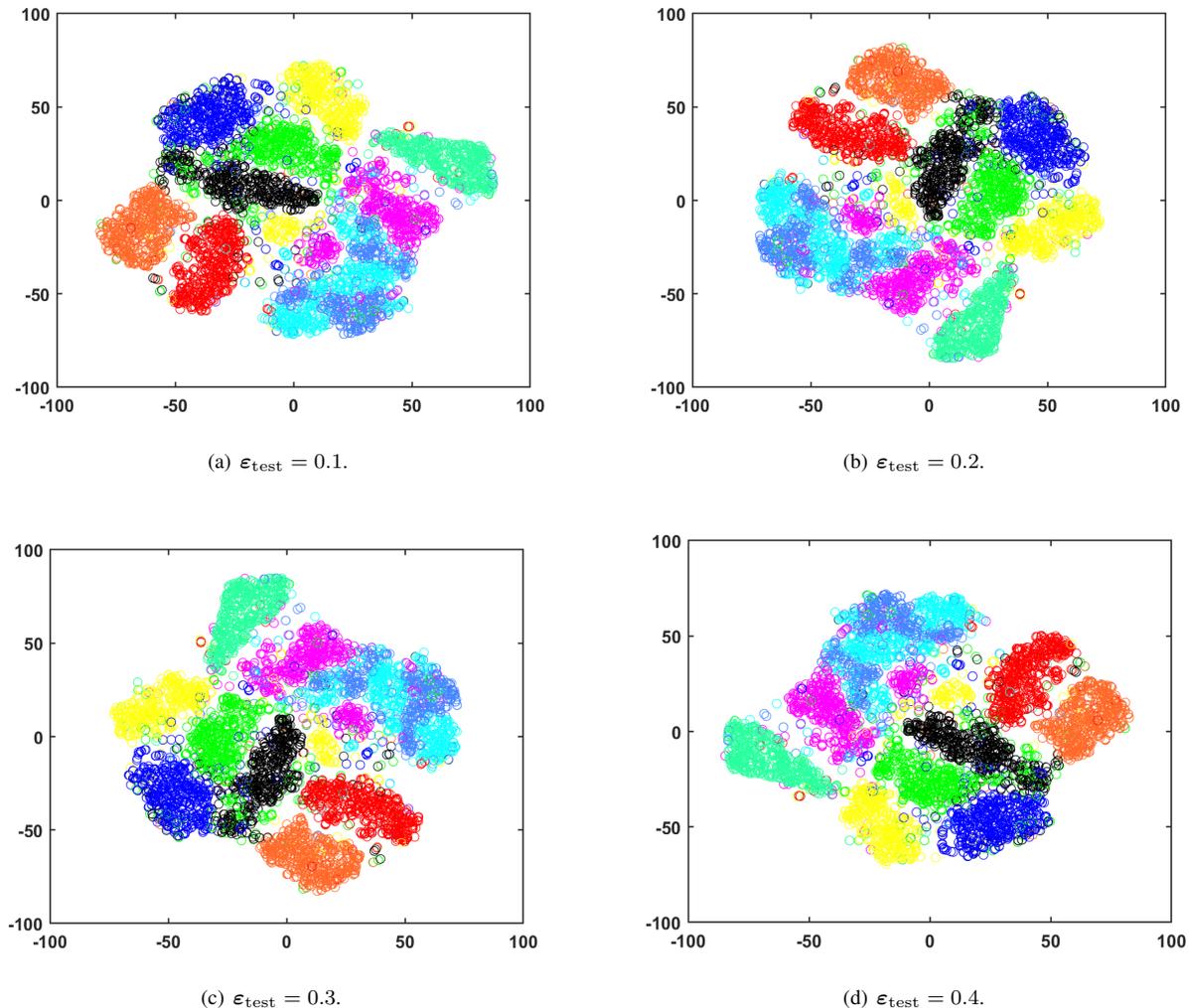
(d) $\varepsilon_{\text{test}} = 0.4$.

Fig. 10. t-SNE visualization of codewords extracted by AIB-JSCC for test set of MNIST dataset. The network is trained with error probability $\varepsilon = 0.1$ and tested with different test error probabilities $\varepsilon_{\text{test}}$. Each color represents a different class.

in MNIST dataset. In particular, we inject noise with different error probabilities into the learned codewords extracted from MNIST and utilize t-Distributed Stochastic Neighbor Embedding (t-SNE) [55] to project the noisy codewords into a 2-dimensional space. From Fig. 10, we can observe that noisy codewords with the same color are close to each other and well separated with other colors, and the distributions of codewords are similar under different $\varepsilon_{\text{test}}$. This is because AIB-JSCC uses the channel conditions in the loss function and enables the learned codewords to resist the channel interference. Therefore, the codewords extracted by AIB-JSCC can preserve semantic information and is robust to different error probabilities.

### B. Parallel-channel Scenario

This subsection evaluates the performance of AIB-JSCC in the parallel-channel scenario. The following three cases with different numbers of subchannels and error probabilities are considered:

$$\varepsilon = \begin{cases} \{0.001, 0.1\} & 2-\text{ch} \\ \{0.001, 0.1, 0.2, 0.25\} & 4-\text{ch} \\ \{0.001, 0.1, 0.2, 0.25, 0.3\} & 5-\text{ch} \end{cases}, \quad (32)$$

where $P - \text{ch}$ represents parallel-channel scenario with $P$ subchannels. Here, the total bandwidth is equally divided into $P$ subchannels.

Table. VII illustrates the reconstruction MSE over test sets. From Table VII, we can observe that AIB-JSCC can achieve better performance in parallel-channel scenarios than that in single-channel scenarios even with smaller error probability. For instance, in the $4 - \text{ch}$ scenario, the average error probability of four subchannels is $0.138$, and the reconstruction error on MNIST and CIFAR are $12.398$ and $52.834$. In contrast, as shown in Table. VII, in the single-channel with smaller error probability, e.g. $\varepsilon = 0.1$, the reconstruction error on MNIST and CIFAR are $12.902$ and $54.464$. This is because in the parallel-channel scenarios, AIB-JSCC utilizes the channel state information in the loss function and is able to transmit important elements over the subchannel with
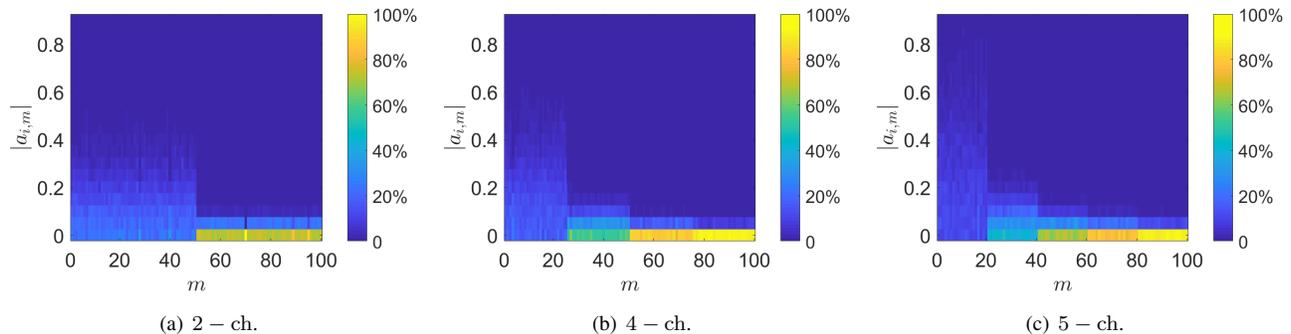
Fig. 11. Distributions of $a_{i,m}$ of different elements in $\hat{\boldsymbol{y}}$ at the decoder under different parallel-channel scenarios.
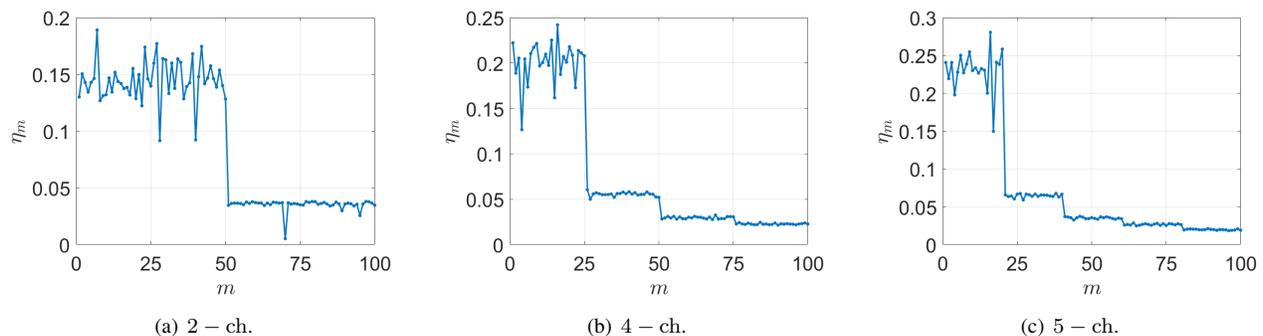


Fig. 12. $\eta_m$ of different elements in $\hat{\boldsymbol{y}}$ at the decoder under different parallel-channel scenarios.

small error probability. Therefore, AIB-JSCC can dynamically allocate elements according to the error probabilities of the subchannels thus improving the reconstruction quality.

Figure 11 shows the distributions of $|a_{i,m}|$, where $a_{i,m}$ represents the weight of the $i$-th neuron of the first FC layer of the decoder at the $m$-th element. From Fig. 11, we can observe that in the subchannel with small error probability, $|a_{i,m}|$ randomly appears in the range of 0 to 0.9. In contrast, in the subchannel with large error probability, $|a_{i,m}|$ concentrates around 0. For example, in Fig. 11(a), in the subchannels with small error probability ($m \leq 50$), $|a_{i,m}|$ mostly appears in the range of 0 to 0.2, and occasionally appears in the range of 0.2 to 0.5. In contrast, in the subchannels with large error probability ($m > 50$), $|a_{i,m}|$ mostly appears in the range of 0 to 0.05 and occasionally appears in the range of 0.05 to 0.1. This implies that the output of the decoder is mainly calculated according to the elements transmitted through subchannel with small error probability, and the elements transmitted through subchannel with large error probability have little effect on the output of the decoder. This is because AIB-JSCC utilizes the channel condition in the loss function and learns to transmit elements important for reconstruction through the subchannel with small error probability to reduce the loss function.

Figure 12 shows the average of $|a_{i,m}|$, $\eta_m$, which is calculated by

$$\eta_m = \frac{1}{L} \sum_{i=1}^{L} |a_{i,m}|, \tag{33}$$

where $L$ represents the number of neurons in the first FC layer of the decoder. From Fig. 12, we can observe that the

elements transmitted through the same subchannel have similar $\eta_m$, and the elements transmitted through the subchannel with small error probability have large $\eta_m$. For example, in Fig. 12(a), in the subchannels with small error probability ($m \leq 50$), $\eta_m$ is in the range of 0.1 to 0.2. In contrast, in the subchannels with large error probability ($m > 50$), $\eta_m$ is in the range of 0 to 0.05, which is much smaller than 0.1. As analyzed before, this is because AIB-JSCC utilizes the channel condition during training and is able to allocate elements important for reconstruction to subchannel with small error probability.

Figure 13 shows the visual reconstructions of AIB-JSCC recovered from noise codewords received from different subchannels. In particular, when using the noisy codeword received from the $i$-th subchannel, i.e., $\hat{\boldsymbol{y}}_{\mathrm{ch}i}$, to reconstruct the images, we fix the other elements in $\hat{\boldsymbol{y}}$ to 0 and feed the new $\hat{\boldsymbol{y}}$ into the trained decoder to obtain the reconstructions. From Fig. 13, we can observe that for both $2 - \mathrm{ch}$ and $5 - \mathrm{ch}$ scenarios, the complete noisy codeword $\hat{\boldsymbol{y}}$ has the best visual performance, and the images recovered from $\hat{\boldsymbol{y}}_{\mathrm{ch}i}$ that is received from the subchannel with smaller error probability, preserve more semantic information. For instance, the images in the second row in Fig. 13(b) can be identified easily, while the images in the third row are difficult to recognize. This is because AIB-JSCC utilizes the channel condition in the loss function, and transmitting the elements with more semantic information through the subchannel with smaller error probability is helpful for semantic information preservation and loss reduction. Therefore, AIB-JSCC is able to transmit the elements with more semantic information for reconstruction
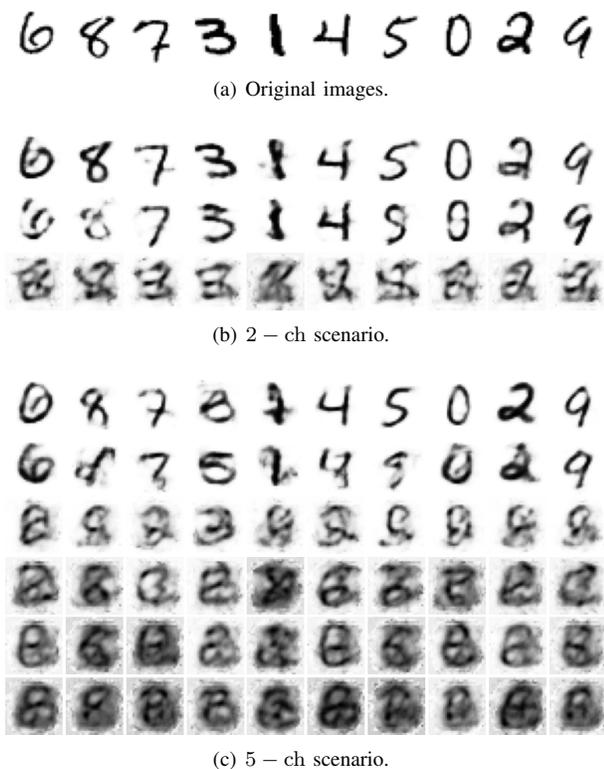
(a) Original images.


(b) $2-$ ch scenario.


(c) $5-$ ch scenario.

Fig. 13. Original image and image recovered in $2-$ ch and $5-$ ch scenarios. For the $2-$ ch and $5-$ ch scenarios, the first row is reconstructed from the complete noisy codeword $\hat{\boldsymbol{y}}$. The second and the third rows in $2-$ ch scenario is reconstructed from $\hat{\boldsymbol{y}}_{\mathrm{ch1}}$ and $\hat{\boldsymbol{y}}_{\mathrm{ch2}}$. The second to the sixth rows in $5-$ ch scenario is reconstructed from $\hat{\boldsymbol{y}}_{\mathrm{ch1}} - \hat{\boldsymbol{y}}_{\mathrm{ch5}}$.

through the subchannel with small error probability.

### C. High-order Scenario

Table VIII shows the MSE of AIB-JSCC and IABF under DMCs with various orders and error probabilities. During training, the orders of the codeword $\boldsymbol{y}$ and the noisy codeword $\hat{\boldsymbol{y}}$ are set to be identical. The channel transition probability is

$$p_{jl} = \begin{cases} 1 - \varepsilon & j = l \\ \frac{\varepsilon}{Q-1} & j \neq l \end{cases}, \tag{34}$$

where $Q$ is the order of $\boldsymbol{y}$. From Table VIII, we can observe that for identical error probability, when the order of $\boldsymbol{y}$ increases, the MSE of AIB-JSCC increases. For instance, when the error probability is $0.1$, the MSE of AIB-JSCC is 57.378, 58.149 and 60.712 when the order of $\boldsymbol{y}$ is 3, 5, 7, respectively. This implies that even though using high-order can improve transmission efficiency, it also diminishes performance. Moreover, AIB-JSCC has a lower MSE than IABF under all considered of DMCs and error probabilities. This is because the proposed IB objective preserves semantic information and discards redundant information.

### VI. CONCLUSION

In this work, we have proposed an AIB-JSCC scheme for image transmission, which can adaptively minimize the transmission rate and distortion at the same time to achieve better reconstruction quality, larger compression ratio, and lower

TABLE VIII
MSE OF RECOVERED IMAGES WHEN THE CHANNEL IS DMC

| Order | Methods | MSE under different error probabilities | | | |
|---|---|---|---|---|---|
| | | 0.1 | 0.2 | 0.3 | 0.4 |
| 3 | IABF | 58.478 | 65.227 | 72.543 | 81.664 |
| | AIB-JSCC | **57.378** | **65.054** | **72.231** | **79.797** |
| 5 | IABF | 59.028 | 66.815 | 74.381 | 81.92 |
| | AIB-JSCC | **58.149** | **65.363** | **73.015** | **80.859** |
| 7 | IABF | 61.901 | 70.097 | 76.169 | 85.379 |
| | AIB-JSCC | **60.712** | **68.463** | **75.777** | **85.175** |

computational complexity than the state-of-the-art approaches. Specifically, we first derived a mathematically tractable form of IB objective for the JSCC system. Then, to appropriately balance the reconstruction distortion and the transmission rate, we further proposed an algorithm that can adaptively adjust hyperparameter $\beta$ of the loss function according to the reconstruction error. Experimental results have shown that with fixed length of codewords, AIB-JSCC always achieved smaller reconstruction error than IB-JSCC and IABF over various error probabilities and datasets, which demonstrates the effectiveness of the proposed IB objective and adaptive IB algorithm. In addition, the images recovered by AIB-JSCC had better visual performance and obtained higher accuracy on downstream classification task than IABF. For a given reconstruction error, AIB-JSCC always permitted larger compression ratio than SSCC and IABF. In particular, AIB-JSCC only needed around $4\%$ and $80\%$ as many elements compared with SSCC and IABF. Moreover, AIB-JSCC also had lower computational complexity and was more robust than IABF. In the parallel-channel scenarios, AIB-JSCC was able to transmit elements important for reconstruction in the subchannel with small error probability. The overall results showed that the proposed schemes can significantly reduce the transmission rate, and improve the reconstruction quality and downstream task accuracy with lower computational complexity.

### APPENDIX
### PROOF OF LEMMA 1

Let $\mathrm{IB}_\beta(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}}) = I(\boldsymbol{x}; \hat{\boldsymbol{y}}) - \beta I(\boldsymbol{x}; \boldsymbol{y})$. We need to guarantee that $\mathrm{IB}_\beta(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})$ is not maximal when $\boldsymbol{x}$ and $\boldsymbol{y}$ are independent, i.e., $p(\boldsymbol{y}|\boldsymbol{x}) = p(\boldsymbol{y})$ or $p(\boldsymbol{x}|\boldsymbol{y}) = p(\boldsymbol{x})$ is not optimal for maximizing $\mathrm{IB}_\beta(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})$. Since there is a Markov chain relationship $\boldsymbol{x} \to \boldsymbol{y} \to \hat{\boldsymbol{y}}$, we have $I(\boldsymbol{x}; \boldsymbol{y}) \geq I(\boldsymbol{x}; \hat{\boldsymbol{y}})$. When $\boldsymbol{x}$ and $\boldsymbol{y}$ are independent, $I(\boldsymbol{x}; \boldsymbol{y}) = I(\boldsymbol{x}; \hat{\boldsymbol{y}}) = 0$ and $\mathrm{IB}_\beta(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})|_{p(\boldsymbol{y}|\boldsymbol{x})=p(\boldsymbol{y})} = 0$ for any $\beta$. Therefore, if $\mathrm{IB}_{\beta_1}(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})$ is not maximal when $p(\boldsymbol{y}|\boldsymbol{x}) = p(\boldsymbol{y})$, there must exist $(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})$ given by $p_1(\boldsymbol{y}|\boldsymbol{x})$ such that

$$\mathrm{IB}_{\beta_1}(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})|_{p(\boldsymbol{y}|\boldsymbol{x})=p_1(\boldsymbol{y}|\boldsymbol{x})} > \mathrm{IB}_\beta(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})|_{p(\boldsymbol{y}|\boldsymbol{x})=p(\boldsymbol{y})} = 0. \tag{35}$$

If $\mathrm{IB}_{\beta_1}(\boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}})$ is not optimal when $p(\boldsymbol{y}|\boldsymbol{x}) = p(\boldsymbol{y})$, we can rewrite (35) as

$$I(\boldsymbol{x}; \hat{\boldsymbol{y}}) - \beta_1 I(\boldsymbol{x}; \boldsymbol{y}) > 0, \exists \boldsymbol{x}, \boldsymbol{y}, \hat{\boldsymbol{y}}. \tag{36}$$

According to (36), we have

$$\beta_1 < \beta_0 = \sup_{\boldsymbol{x} \to \boldsymbol{y} \to \hat{\boldsymbol{y}}} \frac{I\left(\boldsymbol{x}; \hat{\boldsymbol{y}}\right)}{I\left(\boldsymbol{x}; \boldsymbol{y}\right)}. \tag{37}$$

This completes the proof. $\qquad\square$

## References

[1] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Techn. J.*, vol. 27, no. 3, pp. 379–423, Jul. 1948.

[2] E. Berlekamp, R. McEliece, and H. Van Tilborg, "On the inherent intractability of certain coding problems (corresp.)," *IEEE Trans. Inf. Theory*, vol. 24, no. 3, pp. 384–386, May 1978.

[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. 2019.

[4] Y. Yang, C. Guo, F. Liu, C. Liu, L. Sun, Q. Sun, and J. Chen, "Semantic communications with artificial intelligence tasks: Reducing bandwidth requirements and improving artificial intelligence task performance," *IEEE Ind. Electron. Mag.*, to appear, 2022.

[5] M. Chafii, L. Bariah, S. Muhaidat, and M. Debbah, "Ten scientific challenges for 6G: Rethinking the foundations of communications theory," Available: https://arxiv.org/abs/physics/2207.01843, 2022.

[6] A. A. A. Boulogeorgos, J. M. Jornet, and A. Alexiou, "Directional terahertz communication systems for 6G: Fact check?," *IEEE Veh. Technol. Mag.*, vol. 16, no. 4, pp. 68–77, Dec. 2021.

[7] M. Sana and E. C. Strinati, "Learning semantics: An opportunity for effective 6G communications," in *Proc. IEEE Annual Consum. Commun. & Netw. Conf.*, Las Vegas, USA, Feb. 2022, pp. 631–636.

[8] V. Ziegler, H. Viswanathan, H. Flinck, M. Hoffmann, V. Räisänen, and K. Hätönen, "6G architecture to connect the worlds," *IEEE Access*, vol. 8, no. 19981414, pp. 173 508–173 520, Sept. 2020.

[9] Z. Wan, Z. Gao, M. Di Renzo, and L. Hanzo, "The road to industry 4.0 and beyond: A communications-, information-, and operation technology collaboration perspective," Available: https://arxiv.org/abs/physics/2205.04741, 2022.

[10] M. Chen, D. Gündüz, K. Huang, W. Saad, M. Bennis, A. V. Feljan, and H. V. Poor, "Distributed learning in wireless networks: Recent progress and future challenges," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 12, pp. 3579–3605, Oct. 2021.

[11] R. G. Gallager, *Information Theory and Reliable Communication*. Springer, 1968, vol. 2.

[12] M. Gastpar, B. Rimoldi, and M. Vetterli, "To code, or not to code: Lossy source-channel communication revisited," *IEEE Trans. Inf. Theory*, vol. 49, no. 5, pp. 1147–1158, May 2003.

[13] G. Cheung and A. Zakhor, "Bit allocation for joint source/channel coding of scalable video," *IEEE Trans. Image Process.*, vol. 9, no. 3, pp. 340–356, Mar. 2000.

[14] S. Heinen and P. Vary, "Transactions papers source-optimized channel coding for digital transmission channels," *IEEE Trans. Commun.*, vol. 53, no. 4, pp. 592–600, Apr. 2005.

[15] E. Bourtsoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Trans. Cogn. Commun. Netw.*, vol. 5, no. 3, pp. 567–579, Sept. 2019.

[16] D. B. Kurka and D. Gündüz, "DeepJSCC-f: Deep joint source-channel coding of images with feedback," *IEEE J. Sel. Areas Inf. Theory*, vol. 1, no. 1, pp. 178–193, May 2020.

[17] 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.*, vol. 32, no. 4, pp. 2315–2328, Apr. 2021.

[18] K. Choi, K. Tatwawadi, A. Grover, T. Weissman, and S. Ermon, "Neural joint source-channel coding," in *Proc. Int. Conf. Mach. and Learn.*, Long Beach, USA, Jun. 2019, pp. 1182–1192.

[19] Y. Song, M. Xu, L. Yu, H. Zhou, S. Shao, and Y. Yu, "Infomax neural joint source-channel coding via adversarial bit flip," in *Proc. AAAI Conf. Artificial Intell.*, New York, USA, Feb. 2020, pp. 5834–5841.

[20] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Boston, USA, Jun. 2015, pp. 1–9.

[21] D. Minnen, J. Ballé, and G. D. Toderici, "Joint autoregressive and hierarchical priors for learned image compression," in *Neural Inform. Process. Syst.*, vol. 31, Montreal, Canada, Dec. 2018.

[22] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," Available: https://arxiv.org/abs/physics/1810.04805, 2018.

[23] P. Jiang, C. Wen, S. Jin, and G. Y. Li, "Deep source-channel coding for sentence semantic transmission with HARQ," *IEEE Trans. Commun.*, vol. 70, no. 8, pp. 5225–5240, Aug. 2022.

[24] Y. Wang, M. Chen, T. Luo, W. Saad, D. Niyato, H. V. Poor, and S. Cui, "Performance optimization for semantic communications: An attention-based reinforcement learning approach," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 9, pp. 2598–2613, Jul. 2022.

[25] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," Available: https://arxiv.org/abs/physics/0004057, 2000.

[26] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," Available: https://arxiv.org/abs/1612.00410, 2016.

[27] Y. B. Mahabadi, R. Karimi and J. Henderson, "Variational information bottleneck for effective low-resource fine-tuning," Available: https://arxiv.org/abs/2106.05469, 2021.

[28] Y. Du, J. Xu, H. Xiong, Q. Qiu, X. Zhen, C. G. Snoek, and L. Shao, "Learning to learn with variational information bottleneck for domain generalization," in *Proc. Eur. Conf. Comput. Vis.*, Glasgow, UK, Aug. 2020, pp. 200–216.

[29] N. Tishby and N. Zaslavsky, "Deep learning and the information bottleneck principle," in *Proc. IEEE Inf. Theory Worksh.*, Jerusalem, Israel, Apr. 2015, pp. 36–58.

[30] J. Lee, J. Choi, J. Mok, and S. Yoon, "Reducing information bottleneck for weakly supervised semantic segmentation," in *Neural Inform. Process. Syst.*, vol. 34, Virtual, Dec. 2021.

[31] R. Gallager, "Low-density parity-check codes," *IRE Trans. Inf. Theory*, vol. 8, no. 1, pp. 21–28, Jan. 1962.

[32] T. M. Cover, *Elements of Information Theory*. John Wiley & Sons, 1999.

[33] G. Romano and D. Ciuonzo, "Minimum-variance importance-sampling bernoulli estimator for fast simulation of linear block codes over binary symmetric channels," *IEEE Trans. Commun.*, vol. 13, no. 1, pp. 486–496, Dec. 2013.

[34] K. Podgórski, G. Simons, and Y. Ma, "On estimation for a binary-symmetric channel," *IEEE Trans. Inf. Theory*, vol. 44, no. 3, pp. 1260–1272, May 1998.

[35] W. Huleihel and O. Ordentlich, "How to quantize n outputs of a binary symmetric channel to n-1 bits?" in *Proc. Int. Symposium Inf. Theory*, Aachen, Germany, Jun. 2017, pp. 91–95.

[36] D. B. F. Agakov, "The IM algorithm: a variational approach to information maximization," in *Neural Inform. Process. Syst.*, Montreal, Canada, Dec. 2004, p. 201.

[37] A. Mnih and D. Rezende, "Variational inference for monte carlo objectives," in *Proc. Int. Conf. Mach. and Learn.*, New York, USA, Jun. 2016, pp. 2188–2196.

[38] P. Cheng, W. Hao, S. Dai, J. Liu, Z. Gan, and L. Carin, "CLUB: A contrastive log-ratio upper bound of mutual information," in *Proc. Int. Conf. Mach. and Learn.*, Vienna, Austria, Jul. 2020, pp. 1779–1788.

[39] Y. Li, P. Zhao, D. Wang, X. Xian, Y. Liu, and V. S. Sheng, "Learning disentangled user representation based on controllable VAE for recommendation," in *Proc. Int. Conf. Database Syst. Advanced Applicat.*, Taipei, Taiwan, Apr. 2021, pp. 179–194.

[40] T. Wu, I. Fischer, I. L. Chuang, and M. Tegmark, "Learnability for the information bottleneck," in *Proc. Uncertainty in Artificial Intell. Conf.*, Toronto, Canada, Jul. 2020, pp. 1050–1060.

[41] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. Int. Conf. Comput. Statistics*, Paris, France, Aug. 2010, pp. 177–186.

[42] E. Sariyildiz, H. Yu, and K. Ohnishi, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.

[43] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, Dec. 2015.

[44] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *University of Toronto Tech. Rep*, vol. 1, Jan. 2009.

[45] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *Neural Inform. Process. Syst. Worksh.*, Granada, Spain, Dec. 2011, pp. 3730–3738.

[46] G. K. Wallace, "The JPEG still picture compression standard," *IEEE Trans. Consum. Electron.*, vol. 38, no. 1, pp. xviii–xxxiv, Feb. 1992.

[47] M. Rabbani, *Book Review: JPEG2000: Image Compression Fundamentals, Standards and Practice*. SPIE, 2002.

[48] Google, "WebP compression study," Available: https://developers.google.com/speed/webp/docs/webp_study, 2015.

[49] F. Bellard, "BPG image format," Available: https://bellard.org/bpg/, 2018.

[50] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, "Learned image compression with discretized gaussian mixture likelihoods and attention modules," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Seattle, USA, Jun. 2020, pp. 7936–7945.

[51] F. Mentzer, G. D. Toderici, M. Tschannen, and E. Agustsson, "High-fidelity generative image compression," in *Neural Inform. Process. Syst.*, Virtual, Dec. 2020, pp. 11 913–11 924.

[52] M. Lentmaier, A. Sridharan, D. J. Costello, and K. S. Zigangirov, "Iterative decoding threshold analysis for LDPC convolutional codes," *IEEE Trans. Inf. Theory*, vol. 56, no. 10, pp. 5274–5289, Sept. 2010.

[53] J. Hagenauer, E. Offer, and L. Papke, "Iterative decoding of binary block and convolutional codes," *IEEE Trans. Inf. Theory*, vol. 42, no. 2, pp. 429–445, Mar. 1996.

[54] H. Vikalo, B. Hassibi, and T. Kailath, "Iterative decoding for MIMO channels via modified sphere decoding," *IEEE Trans. Commun.*, vol. 3, no. 6, pp. 2299–2311, Spte. 2004.

[55] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE." *J. Mach. Learn. Research*, vol. 9, no. 11, Nov. 2008.

**Lunan Sun** received the B.S. degree from Beijing Jiaotong University, Beijing, China, in 2018. She is currently pursuing the Ph.D. degree with the Beijing Key Laboratory of Network System Architecture and Convergence, Beijing University of Posts and Telecommunications, Beijing, China. Her current research interests include semantic communications, image transmission and deep learning.

**Yang Yang** (Member, IEEE) is currently an Associate Professor with the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications (BUPT), Beijing, China. He received the B.S. degree in information engineering from School of Communication, Xidian University, Xian, China, in June 2013, and the Ph.D degree in Information and Communication Engineering from School of Information and Communication Engineering, Beijing University of Posts and Telecommunications (BUPT), Beijing, China. His research interests include semantic communications, visible light communication and localization. He served as a workshop co-chair/TPC member for a series of IEEE conferences including Globecom, ICC and WCNC. He was a recipient of the IEEE Wireless Communications and Networking Conference (WCNC) 2021 Best Paper Award.

**Mingzhe Chen** (Member, IEEE) is currently an Assistant Professor with the Department of Electrical and Computer Engineering and Institute of Data Science and Computing at University of Miami. His research interests include federated learning, reinforcement learning, virtual reality, unmanned aerial vehicles, and Internet of Things. He has received from the IEEE Communication Society four journal paper awards including the IEEE Marconi Prize Paper Award in Wireless Communications in 2023, the Young Author Best Paper Award in 2021 and 2023, and the Fred W. Ellersick Prize Award in 2022, and three conference best paper awards at IEEE ICC in 2020, IEEE GLOBECOM in 2020, and IEEE WCNC in 2021. He currently serves as an Associate Editor of IEEE Transactions on Mobile Computing, IEEE Wireless Communications Letters, IEEE Transactions on Green Communications and Networking, and IEEE Transactions on Machine Learning in Communications and Networking.

**Caili Guo** (Senior Member, IEEE) received the Ph.D. degree in Communication and Information Systems from Beijing University of Posts and Telecommunication (BUPT) in 2008. She is currently a Professor in the School of Information and Communication Engineering at BUPT. Her general research interests include machine learning and statistical signal processing for wireless communications, with current emphasis on semantic communications, deep learning, and intelligence-enabled edge computing for vehicle communications.

In the related areas, she has published over 200 papers and holds over 30 granted patents. She won Diamond Best Paper Award of IEEE ICME 2018 and Best Paper Award of IEEE WCNC 2021.

**Walid Saad** (Fellow, IEEE) received his Ph.D degree from the University of Oslo, Norway in 2010. He is currently a Professor at the Department of Electrical and Computer Engineering at Virginia Tech, where he leads the Network sciEnce, Wireless, and Security (NEWS) laboratory. He is also the Next-G Wireless Faculty Lead at Virginia Tech's Innovation Campus. His research interests include wireless networks (5G/6G/beyond), machine learning, game theory, security, UAVs, semantic communications, cyber-physical systems, and network science. Dr. Saad is a Fellow of the IEEE. He is also the recipient of the NSF CAREER award in 2013, the AFOSR summer faculty fellowship in 2014, and the Young Investigator Award from the Office of Naval Research (ONR) in 2015. He was the (co-)author of eleven conference best paper awards at IEEE WiOpt in 2009, ICIMP in 2010, IEEE WCNC in 2012, IEEE PIMRC in 2015, IEEE SmartGridComm in 2015, EuCNC in 2017, IEEE GLOBECOM (2018 and 2020), IFIP NTMS in 2019, IEEE ICC (2020 and 2022). He is the recipient of the 2015 and 2022 Fred W. Ellersick Prize from the IEEE Communications Society, and of the IEEE Communications Society Marconi Prize Award in 2023. He was also a co-author of the papers that received the IEEE Communications Society Young Author Best Paper award in 2019, 2021, and 2023. Other recognitions include the 2017 IEEE ComSoc Best Young Professional in Academia award, the 2018 IEEE ComSoc Radio Communications Committee Early Achievement Award, and the 2019 IEEE ComSoc Communication Theory Technical Committee Early Achievement Award. From 2015-2017, Dr. Saad was named the Stephen O. Lane Junior Faculty Fellow at Virginia Tech and, in 2017, he was named College of Engineering Faculty Fellow. He received the Dean's award for Research Excellence from Virginia Tech in 2019. He was also an IEEE Distinguished Lecturer in 2019-2020. He has been annually listed in the Clarivate Web of Science Highly Cited Researcher List since 2019. He currently serves as an Area Editor for the IEEE Transactions on Network Science and Engineering and the IEEE Transactions on Communications. He is the Editor-in-Chief for the IEEE Transactions on Machine Learning in Communications and Networking.

**H. Vincent Poor** (Life Fellow, IEEE) received the Ph.D. degree in EECS from Princeton University in 1977. From 1977 until 1990, he was on the faculty of the University of Illinois at Urbana-Champaign. Since 1990 he has been on the faculty at Princeton, where he is currently the Michael Henry Strater University Professor. During 2006 to 2016, he served as the dean of Princeton's School of Engineering and Applied Science. He has also held visiting appointments at several other universities, including most recently at Berkeley and Cambridge. His research interests are in the areas of information theory, machine learning and network science, and their applications in wireless networks, energy systems and related fields. Among his publications in these areas is the recent book Machine Learning and Wireless Communications. (Cambridge University Press, 2022). Dr. Poor is a member of the National Academy of Engineering and the National Academy of Sciences and is a foreign member of the Chinese Academy of Sciences, the Royal Society, and other national and international academies. He received the IEEE Alexander Graham Bell Medal in 2017.