

Toward Adaptive Semantic Communications: Efficient Data Transmission via Online Learned Nonlinear Transform Source-Channel Coding

Jincheng Dai, *Member, IEEE*, Sixian Wang, *Graduate Student Member, IEEE*, Ke Yang, *Graduate Student Member, IEEE*, Kailin Tan, *Graduate Student Member, IEEE*, Xiaoqi Qin, *Member, IEEE*, Zhongwei Si, *Member, IEEE*, Kai Niu, *Member, IEEE*, and Ping Zhang, *Fellow, IEEE*

Abstract—The emerging field *semantic communication* is driving the research of end-to-end data transmission. By utilizing the powerful representation ability of deep learning models, learned data transmission schemes have exhibited superior performance than the established source and channel coding methods. While, so far, research efforts mainly concentrated on architecture and model improvements toward a static target domain. Despite their successes, such learned models are still suboptimal due to the limitations in model capacity and imperfect optimization and generalization, particularly when the testing data distribution or channel response is different from that adopted for model training, as is likely to be the case in real-world. To tackle this, in this paper, we propose a novel online learned joint source and channel coding approach that leverages the deep learning model’s overfitting property. Specifically, we update the off-the-shelf pre-trained models after deployment in a lightweight online fashion to adapt to the distribution shifts in source data and environment domain. We take the overfitting concept to the extreme, proposing a series of implementation-friendly methods to adapt the codec model or representations to an individual data or channel state instance, which can further lead to substantial gains in terms of the end-to-end rate-distortion performance. Accordingly, the streaming ingredients include both the semantic representations of source data and the online updated decoder model parameters. The system design is formulated as a joint optimization problem whose goal is to minimize the loss function, a tripartite trade-off among the data stream bandwidth cost, model stream bandwidth cost, and end-to-end distortion. The proposed methods enable the communication-efficient adaptation for all parameters in the network without sacrificing decoding speed. Extensive experiments, including user study, on continually changing target source data and wireless channel environments, demonstrate the effectiveness and efficiency of our approach, on which we outperform existing state-of-the-art engineered transmission scheme (VVC combined with 5G LDPC coded transmission).

Index Terms—Semantic communications, online learning, data stream, model stream, end-to-end rate-distortion trade-off.

I. INTRODUCTION

This work was supported in part by the National Natural Science Foundation of China under Grant 62293481, Grant 92067202, Grant 62001049, Grant 62071058, and Grant 61971062, in part by the Beijing Natural Science Foundation under Grant 4222012, in part by Program for Youth Innovative Research Team of BUPT No. 2023QNTD02. (*Corresponding authors: Jincheng Dai, Ping Zhang*)

Jincheng Dai, Sixian Wang, Ke Yang, Kailin Tan, and Zhongwei Si are with the Key Laboratory of Universal Wireless Communications, Ministry of Education, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: daijincheng@bupt.edu.cn).

Xiaoqi Qin, Kai Niu, and Ping Zhang are with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China.

SEMANTIC communications are recently emerging as a new paradigm driving the in-depth fusion of information and communication technology (ICT) advances and artificial intelligence (AI) innovations [1]–[4]. Unlike traditional communication design philosophy that focuses on accurately transmitting bits over a noisy communication channel [5], semantic communications are goal-oriented, which helps the transceiver identify the most valuable information more efficiently, i.e., the information necessary to recover the purpose intended by the transmitter. Performance assessment also goes beyond the common Shannon paradigm of guaranteeing the correct reception of each single transmitted bit, human perceptual loss [6]–[8] and machine task accuracy [9] are taken as the distortion, which is better aligned with the essential goal of end-to-end communications [1].

One roadmap to realize semantic communication is bridging the source and channel parts together to boost the end-to-end content delivery. The paradigm aiming at the integrated design of source and channel processing is *joint source-channel coding (JSCC)* [10], a classical topic in the information theory and coding theory. However, conventional JSCC schemes [10]–[13] are limited by explicit probabilistic models and handcrafted designs, whose optimization is intractable for complex sources. They also ignore the semantic feature aspects of source messages and cannot be optimized towards human perception or machine task directly. In contrast, the emerging semantic communication utilizes deep learning models to realize JSCC [14]–[21], which can be optimized for specific end-to-end transmission objectives. For example, in the case of wireless image transmission, deep JSCC approaches have been verified to surpass classical separation-based BPG source compression [22] combined with advanced low-density parity-check (LDPC) channel coding [23], especially for sources of tiny dimensions, e.g., CIFAR dataset (32×32 pixels). To address the challenge of high resolution media transmission, the optimization goal of JSCC system should be formulated as a trade-off between the reconstruction quality (end-to-end distortion) and the channel bandwidth cost (channel bandwidth ratio). Following this, nonlinear transform source-channel coding (NTSCC) proposed in [24] has achieved content-aware variable-length JSCC via introducing an entropy model on the semantic latent representations, which can significantly improve the overall coding efficiency. This coding paradigm reveals the key aspects of semantic transmission: maximizing reconstruction quality meets the human perception; whereas

minimizing wireless channel bandwidth cost benefits the efficient transmission, i.e., *channel bandwidth ratio-distortion (CBR-D) trade-off optimization*.

Although existing end-to-end transmission approaches have proven to be successful in optimizing the end-to-end *expected CBR-D* trade-off over a source dataset and ergodic wireless channel responses, they are yet unlikely to be optimal for every test instance due to the limited model capacity and imperfect optimization. In essence, they assume that if a model performs well on both training and validation datasets, it will likely generalize well to new, unseen data and channel responses during testing. However, this assumption does not always hold in practice. On the one hand, such an amortized model might not be good at capturing the data semantic feature and channel state for each instance, resulting in suboptimal transform and coding during the inference stage. On the other hand, the imperfect optimization and generalization will be especially severe when the testing data distribution or channel response is different from that adopted in the training stage. To tackle this, we explore a new online learned approach by optimizing the network parameters or the semantic representations during the model inference stage, based on the current target source data and wireless channel domain. In other words, we turn to optimize the CBR-D trade-off for substantial gains on every data and channel state instance.

Some insights in traditional source compression codec can bring some inspirations [25]–[27]. Conventional image/video compressors follow the hybrid transform coding paradigm. For example, HEVC [28] and VVC [29] jointly use discrete cosine transform (DCT) and discrete sine transform (DST) to handle different types of signals. Multiple transform selection (MTS) mechanism is introduced to VVC standard to select the most appropriate transform locally aiming at the best rate-distortion trade-off. Inspired by the idea of signal-dependent transform in traditional source compression methods, neural codecs are also evolving toward the *instance-adaptive* paradigm. The suboptimality of neural codecs has been studied extensively in terms of the model inference suboptimality [30]. It has been shown that by online finetuning the encoder parameters or latent features from a well-trained model for a particular instance, substantial gain can be further obtained in the compression rate-distortion performance [31]–[36]. Some new methods are recently developed to further improve the rate-distortion cost, e.g., a full-model instance-adaptive compression method was proposed in [37], a neural syntax method was proposed in [38] to realize data-dependent compression, etc. All these methods leverage neural network’s overfitting property, adapting the model to an individual source data sample.

Inspired by the useful insights from traditional compression codecs and neural compression codecs, in this paper, we make the first attempt to build a neural instance/domain adaptive joint source-channel coding architecture for end-to-end data transmission. We show that overfitted neural-enhancement on end-to-end communication systems is feasible and effective when combined with online learning. Accordingly, we propose *online learned adaptive NTSCC*, a novel framework that employs online learning to overfit the instant source data sample and channel state information (CSI). Our new method is easy

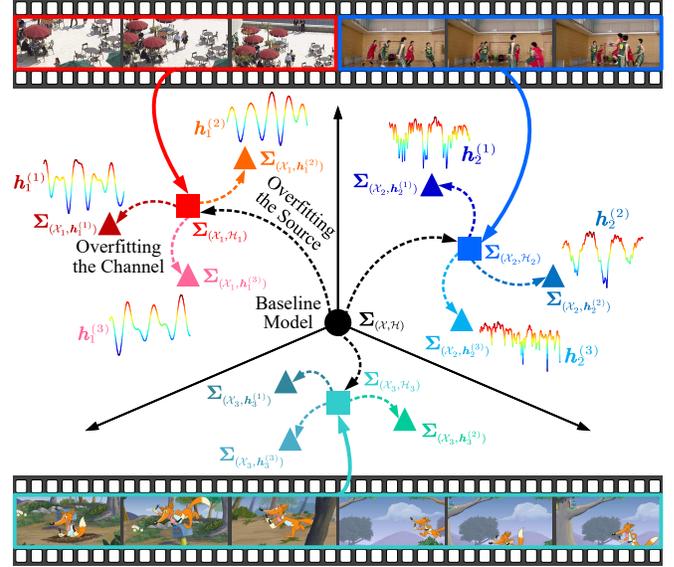


Fig. 1. We first adapt the baseline model $\Sigma(\mathcal{X}, \mathcal{H})$ learned from the source data set \mathcal{X} and CSI set \mathcal{H} to the given specific source domain \mathcal{X}_k and CSI domain \mathcal{H}_k , i.e., overfitted models $\Sigma(\mathcal{X}_k, \mathcal{H}_k)$. To address the problem of transmission over the time-varying wireless channel, we further design a channel-dependent method to adapt the model $\Sigma(\mathcal{X}_k, \mathcal{H}_k)$ to every specific CSI instance $h_k^{(i)} \in \mathcal{H}_k$, resulting in the final overfitted model $\Sigma(\mathcal{X}_k, h_k^{(i)})$.

to implement and can be incorporated in a number of different architectures of semantic communications. In this paper, as a representative case, we inject our online adaptation method to the nonlinear transform source-channel coding (NTSCC) based semantic communication system [24], to validate the effectiveness and efficiency. In contrast to previous works, our system introduces an additional *model stream* alongside the traditional *data stream*, which is utilized to update the JSCC decoder and synthesis transform parameters at the receiver. We take into account the costs of sending the model updates such that the whole system design is formulated as an optimization problem whose goal is to minimize the loss function that is a tripartite trade-off among the data stream bandwidth cost (R), model stream bandwidth cost (M), and end-to-end distortion (D) terms. Our proposed method is lightweight, effective and intuitively appealing, which could be feasibly transplanted to any existing deep learning based semantic communication systems. Fig. 1 illustrates the key idea of our work briefly.

Specifically, the contributions of this paper can be summarized as follows.

- (1) *Online Learned NTSCC Framework*: We make the first attempt to build a source and channel instance or domain adaptive semantic communication system based on the online learned NTSCC. Our innovative overfitting mechanism enables the whole NTSCC system to be much more powerful and flexible for extracting more compact semantic representations, offering superior end-to-end bandwidth ratio-distortion performance. We make variational analysis to interpret the origins of general model inference suboptimality, and clarify the rationale of our proposed source and channel overfitting paradigm.
- (2) *Overfitting the Source*: We design two simple yet efficient

ways to realize the source data instance or domain adaptive NTSCC. These methods are of different ideas for adapting a given pre-trained model or semantic latent representation to an individual data sample or other content domain that is different in appearance. We discuss the specific scenario for which each method is suitable. Our source overfitting mechanism online adapts a baseline model to each specific scene to maximize communication efficiency.

- (3) *Overfitting the Channel:* We propose a plug-in CSI modulation module inserted into pre-trained codec modules. It enables the whole system to efficiently deal with different channel states with a single trained network. The proposed scheme can not only adapt to different signal-to-noise ratio (SNR) under the block fading channel, but also provide consistent and robust performance under diverse frequency selective fading channels, which makes our model agilely transferred over various channel states.
- (4) *Performance Validation:* We verify the effectiveness and efficiency of online learned NTSCC over video I-frame sources and practical wireless channels. Extensive results indicate that our method can lead to substantial gains in the CBR-D performance without sacrificing decoding speed. Equivalently, achieving the same end-to-end transmission performance, the proposed transceiver adaptation scheme can save up to 45% bandwidth cost compared to the state-of-the-art (SOTA) engineered transmission scheme (VVC combined with 5G LDPC coded transmission).

The remainder of this paper is organized as follows. In Section II, we review the architecture and properties of NTSCC system, and analyze its suboptimality using the variational inference. Next, in Section III, we present our source overfitting methods, including different ideas for online overfitting a given pre-trained baseline model to an individual data sample or other content domain. In Section IV, we show our channel overfitting methods, including details on a plugin-in channel modulation module to adapt a pre-trained model to the instant channel state. Section V shows experimental results to quantify our performance gain, and some valuable discussions are also given. Finally, Section VI concludes this paper.

Notational Conventions: Throughout this paper, lowercase letters (e.g., x) denote scalars, bold lowercase letters (e.g., \mathbf{x}) denote vectors. In some cases, x_i denotes the elements of \mathbf{x} , which may also represent a subvector of \mathbf{x} as described in the context. Bold uppercase letters (e.g., \mathbf{X}) denote matrices, and \mathbf{I}_m denotes an m -dimensional identity matrix. $\ln(\cdot)$ denotes the natural logarithm, and $\log(\cdot)$ denotes the logarithm to base 2. p_x denotes a probability density function (pdf) with respect to the random variable x . In addition, $\mathbb{E}(\cdot)$ denotes the statistical expectation operation, and \mathbb{R} denotes the real number set. Finally, $\mathcal{N}(x|\mu, \sigma^2) \triangleq (2\pi\sigma^2)^{-1/2} \exp(-(x-\mu)^2/(2\sigma^2))$ denotes a Gaussian function, and $\mathcal{U}(a-u, a+u)$ stands for a uniform distribution centered on a with the range from $a-u$ to $a+u$.

II. PRELIMINARIES AND MOTIVATION

Built upon the variational auto-encoder (VAE) [39] architecture, NTSCC has shown superior performance on wireless

image and video transmission problems [1], [24], [40]. Owing to the powerful ability of representation learning, NTSCC can well extract the source semantic features and transmit them over the wireless channels efficiently by using variable-length deep JSCC techniques. This method not only achieves comparable or better performance than the SOTA engineered source compression combined with advanced channel coding schemes, but also greatly surpasses plain auto-encoder based deep JSCC methods [14] that directly encode the raw source data rather than its semantic features. In addition, NTSCC shows great potential to achieve lower time complexity due to its efficient parallel computing with deep neural networks (DNNs). The above superiority of NTSCC can well support semantic communications. Therefore, in this paper, we choose NTSCC to build the semantic communication system. Similar to previous works, we take image or video I-frame (intra-coded frame) source as representative, but our work is extensible for other source modalities.

A. NTSCC based Semantic Communication System

The idea of NTSCC stems from the landmark work of Ballé *et al.* on nonlinear transform coding (NTC) [41]–[43]. Given the pristine data sample \mathbf{x} , e.g., an image \mathbf{x} modeled as a vector of pixel intensities $\mathbf{x} \in \mathbb{R}^m$, it is first transformed into semantic latent representation \mathbf{y} using a DNN-based nonlinear analysis transform g_a . In data compression tasks, \mathbf{y} will be quantized as discrete-valued latent representation $\bar{\mathbf{y}}$, followed by entropy encoding [44] to convert $\bar{\mathbf{y}}$ into bit sequence. This bit sequence will be fed into entropy decoding to losslessly reconstruct $\bar{\mathbf{y}}$, and another nonlinear synthesis transform DNN module g_s uses $\bar{\mathbf{y}}$ to reconstruct the decoded data $\hat{\mathbf{x}}$. g_a and g_s are jointly optimized under the rate-distortion constraint. In communication systems, the above source compressive coding paradigm relies heavily on advanced channel coding and signal processing techniques to ensure the transmitted bit sequence to be losslessly recovered. This separation-based approach has been employed in many current communication systems, as the binary representations of various source data can be seamlessly transmitted over arbitrary wireless channels by changing the underlying channel code.

However, with increasing demands on low-latency wireless data delivery applications such as extended reality (XR), the limits of the separation-based design begin to emerge. Current wireless data transmission systems suffer from time-varying channel conditions, in which case the separation-based design leads to significant *cliff-effect* when the channel condition is below the level anticipated by the channel code [14]. Furthermore, the widely-used entropy coding is quite sensitive to the variational estimate of the marginal distribution of the source latent representation. Small perturbations on this marginal can lead to the catastrophic error propagation in entropy decoding [45]. In practice, the small perturbation is often caused by the floating point round-off error [46]. This round-off operation depends heavily on hardware and software platforms, and in various data compression applications, the transceiver may employ different platforms as stated in [46]. As a result, this non-determinism issue in transmitter vs. receiver will lead to severe performance degradation.

To address the above issues, our idea in NTSCC [24] is to *replace quantization and entropy coding* by integrating source coding and channel coding as a trained DNN, resulting in deep JSCC to transmit the latent representation \mathbf{y} directly. We have achieved better end-to-end transmission performance, and the system is robust to unpredictable wireless channels. The whole procedure of NTSCC is depicted in Fig. 2. The latent code \mathbf{y} is fed into both the analysis transform h_a and the deep JSCC encoder f_e . On the one hand, h_a summarizes the distribution of mean values and standard derivations of \mathbf{y} in the hyperprior \mathbf{z} . The transmitter utilizes \mathbf{z} to estimate the mean vector $\boldsymbol{\mu}$ and the standard derivation vector $\boldsymbol{\sigma}$, and use them to determine the bandwidth to transmit the latent representation. On the other hand, f_e encodes \mathbf{y} as the channel-input sequence $\mathbf{s} \in \mathbb{R}^k$, and the received sequence is $\hat{\mathbf{s}} = W(\mathbf{s})$, whose transition probability is $p_{\hat{\mathbf{s}}|\mathbf{s}}(\hat{\mathbf{s}}|\mathbf{s})$. In this paper, we consider the general fading channel model such that the transfer function is $\hat{\mathbf{s}} = W(\mathbf{s}|\mathbf{h}) = \mathbf{h} \odot \mathbf{s} + \mathbf{n}$ where \odot is the element-wise product, \mathbf{h} denotes the CSI vector, and each component of the noise vector \mathbf{n} is independently sampled from a Gaussian distribution, i.e., $\mathbf{n} \sim p_{\mathbf{n}} \triangleq \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$, where σ_n^2 is the noise power. At the receiver, $\hat{\mathbf{s}}$ is further fed into the deep JSCC decoder f_d to reconstruct the latent representation $\hat{\mathbf{y}}$, which is further used by the nonlinear synthesis transform g_s to recover the source data $\hat{\mathbf{x}}$. The whole procedure of NTSCC system is

$$\mathbf{x} \xrightarrow{g_a(\cdot; \phi_g)} \mathbf{y} \xrightarrow{f_e(\cdot; \phi_f)} \mathbf{s} \xrightarrow{W(\cdot|\mathbf{h})} \hat{\mathbf{s}} \xrightarrow{f_d(\cdot; \theta_f)} \hat{\mathbf{y}} \xrightarrow{g_s(\cdot; \theta_g)} \hat{\mathbf{x}}$$

$$\text{with the latent prior } \mathbf{y} \xrightarrow{h_a(\cdot; \phi_h)} \mathbf{z} \xrightarrow{h_s(\cdot; \theta_h)} \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}, \quad (1)$$

where $(\phi, \theta) = (\phi_g, \phi_h, \phi_f, \theta_g, \theta_h, \theta_f)$ encapsulate learnable DNN parameters of each function. The system efficiency is measured by the *channel bandwidth ratio (CBR)* $\rho = k/m$.

The key idea of NTSCC lies in variable-length deep JSCC guided by the latent prior on the semantic feature space. The latent prior $p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z})$ is obtained as

$$p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) = \prod_i \underbrace{\left(\mathcal{N}(y_i | \mu_i, \sigma_i^2) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right)}_{p_{y_i|\mathbf{z}}}(y_i) \quad (2)$$

with $(\boldsymbol{\mu}, \boldsymbol{\sigma}) = h_s(\mathbf{z}; \theta_h)$,

where the convolutional operation “*” with a standard uniform distribution is used to match the prior to the marginal such that the estimated rate $-\log p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z})$ is non-negative. The hyperprior \mathbf{z} is usually transmitted over the digital link as side information due to its small cost, where the quantization $\lfloor \cdot \rfloor$ (rounding to integers) is needed as marked in Fig. 2. A uniformly-noised proxy $\tilde{\mathbf{z}} = \mathbf{z} + \mathbf{o}$ is used to replace the quantized representation $\bar{\mathbf{y}} = \lfloor \mathbf{y} \rfloor$ during model training [43], where o_j is sampled from $\mathcal{U}(-\frac{1}{2}, \frac{1}{2})$. The probability of hyperprior $\tilde{\mathbf{z}}$ is calculated on the fully factorized density $p_{\mathbf{z}} = \prod_j p_{z_j}$ as

$$p_{\mathbf{z}}(\tilde{\mathbf{z}}) = \prod_j \underbrace{\left(p_{z_j|\psi^{(j)}}(z_j|\psi^{(j)}) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right)}_{p_{z_j}}(\tilde{z}_j), \quad (3)$$

where $\psi^{(j)}$ encapsulates all the parameters of $p_{z_j|\psi^{(j)}}$.

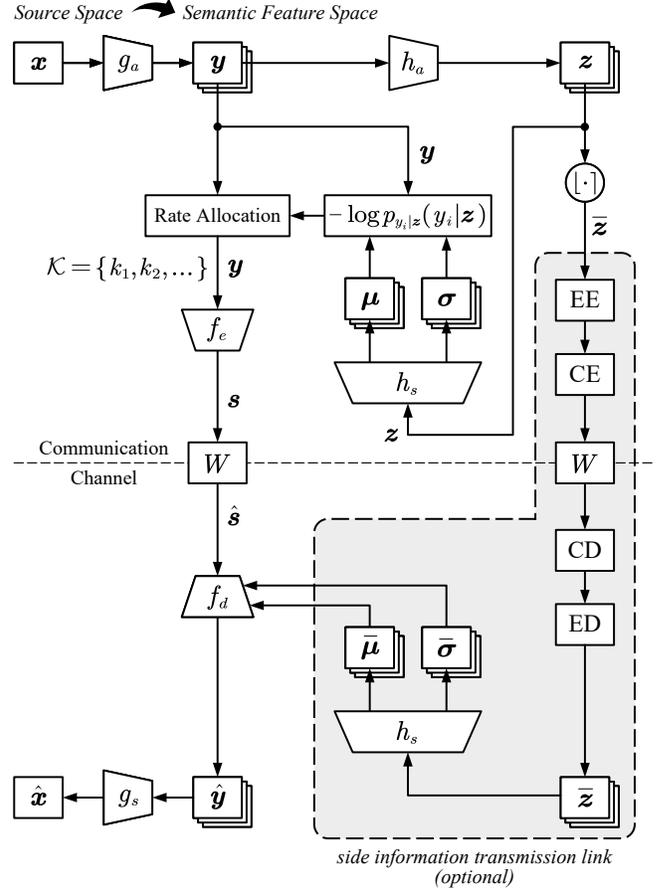


Fig. 2. The whole NTSCC architecture for semantic communications. The transmitter should know the entropy model on $\tilde{\mathbf{z}}$ to entropy encode (EE) and channel encode (CE) it, that is modeled as a non-parametric factorized density conditioned on ψ as (3). The receiver begins with channel decoding (CD) and entropy decoding (ED) to recover the side information $\tilde{\mathbf{z}}$, and then uses it to decode $\hat{\mathbf{y}}$. In addition, the side information $\tilde{\mathbf{z}}$ is not necessary for the receiver as analyzed in [24]. If $\tilde{\mathbf{z}}$ is not transmitted, the decoding performance shows some degradation while the bandwidth cost is also reduced. On the whole, the system bandwidth cost-distortion performance is comparable.

The optimizing problem of NTSCC is formulated following the variational inference context [39], the posterior distribution $p_{\hat{\mathbf{s}}, \tilde{\mathbf{z}}|\mathbf{x}}$ is approximated using the variational density $q_{\hat{\mathbf{s}}, \tilde{\mathbf{z}}|\mathbf{x}}$ by minimizing their Kullback-Leibler (KL) divergence over the data distribution $p_{\mathbf{x}}$ and the CSI distribution $p_{\mathbf{h}}$ as the equation (11) in [24]. Accordingly, the optimization of NTSCC system can be formally converted to the minimization of the expected channel bandwidth cost, as well as the expected distortion of the reconstructed data versus the original, which leads to the optimization of the following R-D trade-off,

$$\mathcal{L}_{\text{R-D}}(\phi, \theta, \psi) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \mathbb{E}_{\mathbf{h} \sim p_{\mathbf{h}}} D_{\text{KL}}(q_{\hat{\mathbf{s}}, \tilde{\mathbf{z}}|\mathbf{x}} \| p_{\hat{\mathbf{s}}, \tilde{\mathbf{z}}|\mathbf{x}}) \Leftrightarrow \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \mathbb{E}_{\mathbf{h} \sim p_{\mathbf{h}}} \left(\underbrace{\lambda \left(-\eta_y \log p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) - \eta_z \log p_{\mathbf{z}}(\tilde{\mathbf{z}}) \right)}_{\text{data stream channel bandwidth cost: } R} + \underbrace{d(\mathbf{x}, \hat{\mathbf{x}})}_{\text{distortion: } D} \right), \quad (4)$$

where the Lagrange multiplier λ on the total channel bandwidth cost determines the trade-off between the data stream bandwidth cost R and the end-to-end distortion D . The scaling factors η_y and η_z control the relation between the estimated entropy and the allocated channel bandwidth, which are tied with

the source-channel codec capability and the wireless channel state. A larger η_y indicates a better performance on deep JSCC codec f_e and f_d , but incurs more channel bandwidth cost. Accordingly, η_y can be adjusted as a hyperparameter to control the system R-D trade-off. η_z is not adjusted manually since explicit entropy coding and LDPC coding are selected to transmit the side information.

In practice, each embedding y_i is a c -dimensional feature vector. The learned entropy model $-\log p_{y_i|z}(y_i|z)$ indicates the summation of entropy along c dimensions of y_i , thus, the information density distribution of \mathbf{y} is captured. Accordingly, the bandwidth cost, such as the number of OFDM subcarriers, \bar{k}_i for transmitting y_i can be determined as

$$\bar{k}_i = Q(k_i) = Q\left(\underbrace{-\eta_y \log p_{y_i|z}(y_i|z)}_{k_i}\right), \quad (5)$$

where the learned entropy model $p_{y_i|z}$ follows (2), Q denotes a scalar quantization whose range includes 2^q ($q = 1, 2, \dots$) integers, and the quantization value set $\mathcal{V} = \{v_1, v_2, \dots, v_{2^q}\}$ is related to the scaling factor η_y and the Lagrange multiplier λ . Hence, the predetermined q bits should be transmitted as extra side information to inform the receiver which bandwidth is allocated to every embedding y_i . To adaptively map y_i to a \bar{k}_i -dimensional channel-input vector s_i , the dynamic neural network structure [47] is introduced into Transformers [48] to realize the deep JSCC codec f_e and f_d [24].

B. Motivation of Online Learned Adaptive NTSCC

In existing works, the loss function \mathcal{L}_{R-D} in (4) is optimized over a corpus of source data samples (such as a large amount of images) and channel states in order to find optimal codec function parameters $\Sigma = (\phi, \theta, \psi)$. Although the models have been trained over a large corpus of source and channel samples for finding out what should be ideally optimal codec functions $(g_a, g_s, h_a, h_s, f_e, f_d)$ over the whole data set and ergodic channel responses, we will show the codec functions can still be improved for each single data sample and instant CSI. This suboptimality of NTSCC model can be interpreted from the *inference suboptimality* of VAE [30] as follows: the mismatch between the true and approximate posterior. It has been proven that the inference gap includes two components: the *approximation gap* and the *amortization gap*. This approximation gap comes from the inability of the variational distribution family to exactly match the true posterior, and the amortization gap refers to the difference caused by amortizing the variational parameters over the entire training set, instead of optimizing for each training example individually.

Specifically, in our NTSCC system, given the source data sample \mathbf{x} and the instant CSI vector \mathbf{h} , the inference gap \mathcal{G} is

$$\mathcal{G} = D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}), \quad (6)$$

where $q_{\hat{s}, \tilde{z}|\mathbf{x}}$ is derived using the parameters (ϕ, θ, ψ) learned under the entire training set, e.g.,

$$q_{\hat{s}, \tilde{z}|\mathbf{x}} = \arg \min_{q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \mathbb{E}_{\mathbf{h} \sim p_{\mathbf{h}}} D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}). \quad (7)$$

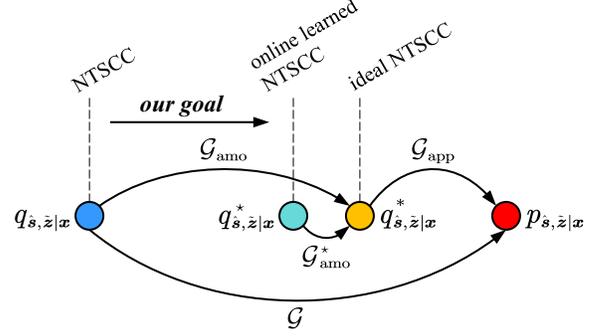


Fig. 3. The relation of gaps in model inference. By overfitting the source and channel, we can upgrade the standard NTSCC system to the online learned NTSCC system, thus making $\mathcal{G}_{\text{amo}}^* \ll \mathcal{G}_{\text{amo}}$.

However, the optimal $q_{\hat{s}, \tilde{z}|\mathbf{x}}^*$ should be derived under the given \mathbf{x} and \mathbf{h} as

$$q_{\hat{s}, \tilde{z}|\mathbf{x}}^* = \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}), \quad (8)$$

which corresponds to the optimal parameters $(\phi^*, \theta^*, \psi^*)$. As illustrated in Fig. 3, it can be derived by

$$\begin{aligned} \mathcal{G} &= D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}) = \underbrace{D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}}^* \| p_{\hat{s}, \tilde{z}|\mathbf{x}})}_{\text{approximation gap: } \mathcal{G}_{\text{app}}} + \\ &\underbrace{D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}) - D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}}^* \| p_{\hat{s}, \tilde{z}|\mathbf{x}})}_{\text{amortization gap: } \mathcal{G}_{\text{amo}}} = \mathcal{G}_{\text{app}} + \mathcal{G}_{\text{amo}}. \end{aligned} \quad (9)$$

This indicates that the amortized posterior q over the entire source dataset and ergodic CSI still incurs the amortization gap \mathcal{G}_{amo} under the instant \mathbf{x} and \mathbf{h} . The above analysis demonstrates that a neural-network-based semantic communication model is trained on the entire dataset and CSI set with the target of achieving the best end-to-end R-D performance on test data and CSI, i.e., we ideally expect $\mathcal{G}_{\text{amo}} = 0 \Leftrightarrow \mathcal{G} = \mathcal{G}_{\text{app}}$. However, due to limited model capacity, optimization difficulties, and insufficient data and CSI, the model cannot in general achieve this goal. When the source data distribution or channel model differs from that in the training phase, model generalization will not be guaranteed even with the infinite data and model capacity, and perfect optimization.

We however note that a convenient feature of neural wireless data transmission is a model or semantic latent representation that can be easily finetuned on new data and CSI. A model can for instance or domain be trained after deployment. Inspired by this, as illustrated in Fig. 3, our goal is closing the amortization gap \mathcal{G}_{amo} by overfitting the source and channel such that the resulting posterior $q_{\hat{s}, \tilde{z}|\mathbf{x}}^*$ of our online learned NTSCC model can approach the optimal $q_{\hat{s}, \tilde{z}|\mathbf{x}}^*$ at every test instance. Thus the new amortization gap $\mathcal{G}_{\text{amo}}^*$ is much smaller than \mathcal{G}_{amo} . As such, we are effectively trying to solve the following optimization problem, given the instant \mathbf{x} and \mathbf{h} ,

$$\begin{aligned} (\phi^*, \theta^*, \psi^*) &= \arg \min_{\phi, \theta, \psi} D_{\text{KL}}(q_{\hat{s}, \tilde{z}|\mathbf{x}} \| p_{\hat{s}, \tilde{z}|\mathbf{x}}) \\ &= \arg \min_{\phi, \theta, \psi} \mathcal{L}_{R-D}(\phi, \theta, \psi, \mathbf{x}, \mathbf{h}) = \arg \min_{\phi, \theta, \psi} \\ &\left(\underbrace{\lambda \left(-\eta_y \log p_{\mathbf{y}|z}(\mathbf{y}|z) - \eta_z \log p_z(\tilde{z}) \right)}_{\text{data stream bandwidth cost: } R} + \underbrace{d(\mathbf{x}, \hat{\mathbf{x}})}_{\text{distortion: } D} \right), \end{aligned} \quad (10)$$

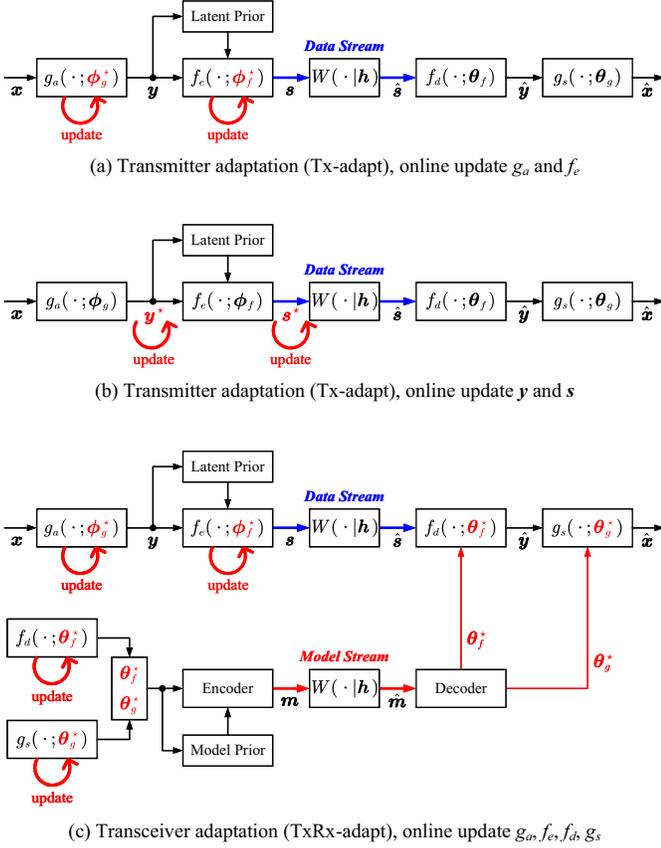


Fig. 4. Key ideas of overfitting the source in the online learned NTSCC.

where $\mathcal{L}_{R-D}(\phi, \theta, \psi, x, h)$ is the transmission R-D objective for a particular x and h .

III. OVERFITTING THE SOURCE

In this section, we present two methods to overfit a single source data instance x (e.g. an image) or data from a specific domain \mathcal{X}' (e.g. a set of I-frames from video sequences in the same scene). They stem from two different ideas:

- *Transmitter adaptation (Tx-adapt):* Given an instance x or a group of domain samples from \mathcal{X}' , online update (g_a, f_e) or (y, s) using gradient descent based on the off-the-shelf pre-trained model.
- *Transceiver adaptation (TxRx-adapt):* Given a group of domain samples from \mathcal{X}' , online update (g_a, f_e, f_d, g_s) using gradient descent based on the off-the-shelf pre-trained model.

Fig. 4 presents the key idea of our proposed methods, where red lines mark the updated ingredients during model inference phase. In the following subsections, the instant CSI h will be sampled from a specific channel model \mathcal{H}' , we present how to overfit the source data instance or a specific domain.

Our goal is adapting the baseline model $\Sigma_{(\mathcal{X}, \mathcal{H})}$ learned on the training dataset \mathcal{X} and channel state domain \mathcal{H} to the given specific source domain \mathcal{X}' and channel domain \mathcal{H}' . In general, \mathcal{X}' can be a subset of \mathcal{X} or a related set to \mathcal{X} , and

\mathcal{H}' is a partial scene of \mathcal{H} , e.g., the CSI set of a local area. In this section, we focus on solving the following problems,

$$\text{source instance adaptation: } \Sigma_{(\mathcal{X}, \mathcal{H})} \Rightarrow \Sigma_{(x, \mathcal{H}')} \quad (11a)$$

$$\text{source domain adaptation: } \Sigma_{(\mathcal{X}, \mathcal{H})} \Rightarrow \Sigma_{(\mathcal{X}', \mathcal{H}')} \quad (11b)$$

Our models $\Sigma_{(x, \mathcal{H}')}$ and $\Sigma_{(\mathcal{X}', \mathcal{H}')}$ are both incidentally adapted to a local channel domain \mathcal{H}' since we assume the exact CSI is only available at the receiver (CSIR) instead of the transmitter. In this case, the transmitter cannot obtain the exact CSI instance $h \in \mathcal{H}'$ in real time. Thus, in our source overfitting algorithms described following, the CSI vector h is sampled from the specific channel domain \mathcal{H}' known at the transmitter. As for how to adapt our model to the exact CSI vector, it will be discussed in the next section.

A. Transmitter Adaptation

Given the source sample x and the baseline model (ϕ, θ, ψ) learned on the entire dataset, our goal is adapting the transmitter model parameters (ϕ_g, ϕ_f) or the latent code and channel-input codeword (y, s) to every single data sample.

The key benefit of transmitter adaptation is achieving an improved end-to-end transmission performance while keeping the predictive model fixed such that the computing time at the receiver stays unchanged. The refined nonlinear analysis transform g_a or the latent code y provides a more compact semantic representation of the source sample x , leading to the superior end-to-end R-D performance. As a trade-off, the transmitter incurs additional encoding time and computational expense, while the decoding delay stays unchanged.

This adaptation mode is suitable for situations where sufficient transmitter side computational power is available (e.g., the media server, often located in a relatively well-provisioned facility, such as the cloud), and the encoding delay is not a pressing issue. As a typical case, existing high-resolution image/video wireless delivery depends critically on the bandwidth resource [49], and video content is known in advance at the capable server side. In this case, the transmitter adaptation paradigm can significantly enhance user quality of experience by utilizing server computation.

For model adaptation, we aim at effectively solving the following optimization problem. During the model inference time, for a single data instance x :

$$(\phi_g^*, \phi_f^*) = \arg \min_{\phi_g, \phi_f} \mathcal{L}_{R-D}(\phi, \theta, \psi, x, h), \quad (12)$$

where we adapt g_a and f_e while fixing the entropy model h_a and h_s . This design aims to reduce the model updating complexity. In addition, if the side information \bar{z} is transmitted, our method ensures that no model updates have to be transmitted to update h_s at the receiver. Experiments can verify that our simplified method achieves comparable performance as that of finetuning the whole transmitter model $\phi = (\phi_g, \phi_f, \phi_h)$. In this work, we solve this problem (12) in an iterative procedure as that in [32]. We apply gradient descent on ϕ_g and ϕ_f to update the nonlinear analysis transform g_a and the deep JSCC encoder f_e . The iterative parameter updating procedures are

$$\phi_g^{(t)} = \phi_g^{(t-1)} - \gamma \nabla_{\phi_g} \mathcal{L}_{R-D}(\phi, \theta, \psi, x, h), \quad (13a)$$

Algorithm 1: Tx Model Online Adaptation

Input: Baseline model parameters (ϕ, θ, ψ) trained on training set, the data sample to be transmitted \mathbf{x} , the CSI vector \mathbf{h} .

Output: The updated models ϕ_g^* and ϕ_f^* .

```

1 procedure ONLINEUPDATEMODELS( $\phi, \theta, \psi, \mathbf{x}, \mathbf{h}$ )
2   Initialize model parameters:  $\phi_g^{(0)} \leftarrow \phi_g$ ,
    $\phi_f^{(0)} \leftarrow \phi_f$ ;
3   for  $t = 1, 2, \dots, T_{\max}$  do
4     Forward pass:  $\mathbf{x} \xrightarrow{g_a(\cdot; \phi_g^{(t-1)})} \mathbf{y} \xrightarrow{f_e(\cdot; \phi_f^{(t-1)})} \hat{\mathbf{s}}$ 
    $\mathbf{s} \xrightarrow{W(\cdot; \mathbf{h})} \hat{\mathbf{s}} \xrightarrow{f_d(\cdot; \theta_f)} \hat{\mathbf{y}} \xrightarrow{g_s(\cdot; \theta_g)} \hat{\mathbf{x}}$  with the
   latent prior  $\mathbf{y} \xrightarrow{h_a(\cdot; \phi_h)} \mathbf{z} \xrightarrow{h_s(\cdot; \theta_h)} \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}$ ;
5     Compute loss  $\mathcal{L}_{\text{R-D}}$  according to (10);
6     Update  $\phi_g^{(t)}$  and  $\phi_f^{(t)}$  using gradients  $\nabla_{\phi_g} \mathcal{L}_{\text{R-D}}$ 
   and  $\nabla_{\phi_f} \mathcal{L}_{\text{R-D}}$  as (13);
7   return Updated models  $\phi_g^* \leftarrow \phi_g^{(T_{\max})}$  and
    $\phi_f^* \leftarrow \phi_f^{(T_{\max})}$ .

```

$$\phi_f^{(t)} = \phi_f^{(t-1)} - \gamma \nabla_{\phi_f} \mathcal{L}_{\text{R-D}}(\phi, \theta, \psi, \mathbf{x}, \mathbf{h}), \quad (13b)$$

where γ denotes the learning rate. The final pipeline of Tx model adaptation is described in Algorithm 1.

Apparently, instance adaptation is an extreme case of domain adaptation, which provides better performance with the expense of more encoding time. In practice, we pursue for more efficient online adaptation, thus the Tx model adaptation in Algorithm 1 is often invoked only for several spare source samples from the same domain, e.g., a set of I-frames from a single video whose content locates in the same scene. We also note that since the decoder parameters remain unchanged such that the performance improvements are limited. A more effective transceiver adaptation mode tailored for domain adaptation will be explored in the subsequent subsection.

Next, we introduce a more lightweight method to adapt the latent representation \mathbf{y} and the deep JSCC codeword \mathbf{s} for instance adaptation only. This code adaptation method aims to find more compact semantic representations directly without changing the encoder/decoder parameters. We apply gradient descent on \mathbf{y} and \mathbf{s} to update the semantic latent representation and the deep JSCC codeword, respectively. A special note is that the iterative updating procedure on these two terms should be executed sequentially because \mathbf{s} is generated from \mathbf{y} by using the function f_e . The iterative updating procedure can be written as

$$\mathbf{y}^{(t)} = \mathbf{y}^{(t-1)} - \gamma \nabla_{\mathbf{y}} \mathcal{L}_{\text{R-D}}(\phi, \theta, \psi, \mathbf{x}, \mathbf{h}), \quad (14)$$

with $t = 1, 2, \dots, Y_{\max}$.

After t reaches Y_{\max} , the procedure turns to update \mathbf{s} , i.e.,

$$\mathbf{s}^{(t)} = \mathbf{s}^{(t-1)} - \gamma \nabla_{\mathbf{s}} \mathcal{L}_{\text{R-D}}(\phi, \theta, \psi, \mathbf{x}, \mathbf{h}), \quad (15)$$

with $t = 1, 2, \dots, S_{\max}$.

The total number of updating steps is $T_{\max} = Y_{\max} + S_{\max}$. The final pipeline of Tx code adaptation is described by Al-

Algorithm 2: Tx Code Online Adaptation

Input: Baseline model parameters (ϕ, θ, ψ) trained on training set, the data sample to be transmitted \mathbf{x} , the CSI vector \mathbf{h} .

Output: The updated deep JSCC codeword \mathbf{s}^* .

```

1 procedure ONLINEUPDATELATENTS( $\phi, \theta, \psi, \mathbf{x}, \mathbf{h}, \mathbf{y}$ )
2   Initialize latent representation:  $\mathbf{y}^{(0)} \leftarrow \mathbf{y}$ ;
3   for  $t = 1, 2, \dots, Y_{\max}$  do
4     Forward pass:  $\mathbf{y}^{(t-1)} \xrightarrow{f_e(\cdot; \phi_f)} \mathbf{s} \xrightarrow{W(\cdot; \mathbf{h})} \hat{\mathbf{s}}$ 
    $\hat{\mathbf{s}} \xrightarrow{f_d(\cdot; \theta_f)} \hat{\mathbf{y}} \xrightarrow{g_s(\cdot; \theta_g)} \hat{\mathbf{x}}$  with the latent prior
    $\mathbf{y}^{(t-1)} \xrightarrow{h_a(\cdot; \phi_h)} \mathbf{z} \xrightarrow{h_s(\cdot; \theta_h)} \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}$ ;
5     Compute loss  $\mathcal{L}_{\text{R-D}}$  according to (10);
6     Update  $\mathbf{y}^{(t)}$  using gradients  $\nabla_{\mathbf{y}} \mathcal{L}_{\text{R-D}}$  as (14);
7   return Updated latent representation
    $\mathbf{y}^* \leftarrow \mathbf{y}^{(Y_{\max})}$ .
8 procedure ONLINEUPDATECODE-
   WORDS( $\phi, \theta, \psi, \mathbf{x}, \mathbf{h}, \mathbf{s}$ )
9   Initialize deep JSCC codeword:  $\mathbf{s}^{(0)} \leftarrow \mathbf{s}$ ;
10  for  $t = 1, 2, \dots, S_{\max}$  do
11    Forward pass:
    $\mathbf{s}^{(t-1)} \xrightarrow{W(\cdot; \mathbf{h})} \hat{\mathbf{s}} \xrightarrow{f_d(\cdot; \theta_f)} \hat{\mathbf{y}} \xrightarrow{g_s(\cdot; \theta_g)} \hat{\mathbf{x}}$ ;
12    Compute loss  $\mathcal{L}_{\text{R-D}}$  according to (10);
13    Update  $\mathbf{s}^{(t)}$  using gradients  $\nabla_{\mathbf{s}} \mathcal{L}_{\text{R-D}}$  as (15);
14  return Updated deep JSCC codeword
    $\mathbf{s}^* \leftarrow \mathbf{s}^{(S_{\max})}$ .

```

gorithm 2. This latent representation and codeword adaptation technique reduces the number of updated parameters for lower complexity and decreased GPU peak memory, but it is only applicable for instance adaptation, not for domain adaptation.

B. Transceiver Adaptation

The above transmitter adaptation method is appealing since no additional information needs to be added to the wireless transmitted signal, and nothing changes on the receiver. However, performance gains are relatively limited since the deep JSCC decoder f_d and the nonlinear synthesis transform g_s cannot be adapted. In this subsection, we present a method for transceiver *full-model adaptation*, which tailors the entire NTSCC model to a specific domain. Unlike previous methods, our adaptive NTSCC with the full-model adaptation involves both *data stream* and *model stream* as the wireless transmitted signal. The model stream is utilized to inform the receiving end to update the model parameters of f_d and g_s . *A noteworthy point is that the transceiver full-model adaptation method is particularly suited for domain adaptation, as the model stream bandwidth cost associated with instance adaptation cannot be effectively amortized and results in inefficient computational complexity and high model stream bandwidth cost.*

In this manner, the adapted model can be applied to other unseen samples within the same domain. Also, the full-model online learning process only takes place at the transmitter, as we assume the transmitter has a local copy of the pre-trained

decoder models. This adaptation paradigm yields two kinds of tradeoffs:

- 1) *R-D-M trade-off*: a tripartite trade-off among the averaged data stream bandwidth cost (R), model stream bandwidth cost per domain (M), and distortion (D) terms, formulating the end-to-end R-D-M loss. Consider a specific source domain including N samples, the model stream is transmitted only once, resulting in the actually averaged channel bandwidth cost of $k = R + M/N$, such that the averaged CBR is $\rho = k/m = (R + M/N)/m$.
- 2) *Performance-complexity trade-off*: the more constrained the domain adaptation, the greater potential gains from adaptation. However, a more restrictive domain necessitates multiple updating processes at the transmitter and multiple transmission model updates, leading to increased encoding complexity and delay.

Transceiver full-model adaptation online updates a set of global baseline model parameters (ϕ, θ, ψ) on a single source data instance \mathbf{x} . In practice, similar to that in the transmitter adaptation, we also fix the entropy model h_a and h_s to reduce the adaptation complexity. This results in the updated parameters $(\phi_g^*, \phi_f^*, \theta_g^*, \theta_f^*)$, of which only θ_g^* and θ_f^* are transmitted over the wireless channel as the model stream.

Intuitively, a larger model steam bandwidth cost M offers greater flexibility during model updating, potentially reducing the amortization gap for improved overall R-D performance. However, the actual performance, accounting for the model stream transmission cost, still requires further evaluation. As a result, developing an efficient model stream transmission method along is both challenging and crucial. Inspired by the residual coding idea, we transmit only the changes relative to the baseline model $\delta_g = \theta_g^* - \theta_g$ and $\delta_f = \theta_f^* - \theta_f$ in practice. To encode the model updates $\delta = (\delta_g, \delta_f)$, we need to build a model prior $p_\delta(\delta)$ to quantify the model rate. Accordingly, the model rate is derived with $-\log p_\delta(\delta)$. Adding this term to the R-D loss function in (10), we obtain the full-model adaptive NTSCC objective:

$$\begin{aligned} \mathcal{L}_{\text{R-D-M}}(\phi, \theta, \psi, \delta, \mathbf{x}, \mathbf{h}) &= \mathcal{L}_{\text{R-D}}(\phi, \theta + \hat{\delta}, \psi, \mathbf{x}, \mathbf{h}) + \beta(-\eta_\delta \log p_\delta(\delta)) \\ &= \lambda \underbrace{(-\eta_y \log p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) - \eta_z \log p_z(\tilde{\mathbf{z}}))}_{\text{data stream bandwidth cost: } R} + \underbrace{d(\mathbf{x}, \hat{\mathbf{x}})}_{\text{distortion: } D} \\ &\quad + \beta \underbrace{(-\eta_\delta \log p_\delta(\delta))}_{\text{model stream bandwidth cost: } M}, \end{aligned} \quad (16)$$

where the scaling factor η_δ is tied with the capability of codec used to transmit δ , $\hat{\delta}$ denotes the reconstructed δ at the receiver end, and β controls the trade-off between the standard R-D loss and the model stream bandwidth cost. Minimization of $\mathcal{L}_{\text{R-D-M}}$ in (16) ensures any cost in the model stream contributes to the R-D performance improvement.

For the model prior p_δ , any probability distribution function can be selected, herein, we naturally define p_δ as a factorized model, i.e., $p_\delta(\delta) = \prod_i p_{\delta_i}(\delta_i)$, where every component is of the same parameter. Each $p_{\delta_i}(\delta_i)$ is consistently generated

from zero-centered Gaussian with a shared variance σ^2 as

$$\begin{aligned} p_\delta(\delta) &= \prod_i \underbrace{\left(\mathcal{N}(\delta_i|0, \sigma^2) * \left(\Delta \cdot \mathcal{U}\left(-\frac{\Delta}{2}, \frac{\Delta}{2}\right) \right) \right)}_{p_{\delta_i}}(\delta_i) \\ &= \prod_i \int_{\delta_i - \frac{\Delta}{2}}^{\delta_i + \frac{\Delta}{2}} \mathcal{N}(\delta'_i|0, \sigma^2) d\delta'_i, \end{aligned} \quad (17)$$

where the uniform distribution convolution is utilized to relax the prior such that the estimated model rate $-\log p_\delta(\delta)$ stays non-negative. It can directly interpolate the discrete probability values $p_{\delta_i}(\delta_i)$ at the quantized values $\bar{\delta}_i$ that will be used when δ is transmitted over the digital link with entropy coding and channel coding, and Δ in (17) indicates the quantization bin width. If δ needs to be quantized for entropy coding, we define the quantization function as that in [37] with N width- Δ quantization bins, i.e.,

$$\begin{aligned} \bar{\delta}_i &= Q_\Delta(\delta_i) \\ &= \text{clip}\left(\left\lfloor \frac{\delta_i}{\Delta} \right\rfloor \cdot \Delta, \min = -\frac{(N-1)\Delta}{2}, \max = \frac{(N-1)\Delta}{2}\right). \end{aligned} \quad (18)$$

During model training, the gradient of Q_Δ can be approximated by the Straight-Through estimator (STE) [50]. Correspondingly, we also leverage the uniformly-noised proxy $\tilde{\delta} = \delta + \mathbf{o}$ with $o_i \sim \mathcal{U}(-\frac{\Delta}{2}, \frac{\Delta}{2})$ to compute the model stream rate during online training. In this case, the model stream rate is evaluated by substituting $\tilde{\delta}$ into the model prior p_δ as $-\log p_\delta(\tilde{\delta})$ in the R-D-M loss (16).

Note that most model updates δ_i (or the quantized $\bar{\delta}_i$) center around zero, the zero-centered Gaussian model prior in (17) indeed ensures the minimum model stream cost for all-zero update, i.e., $-\log p_\delta(\mathbf{0})$. To reduce the model stream cost, we can further generalize the standard Gaussian model prior as a *Gaussian mixed model (GMM)* prior, which stems from [37]. In particular, we adopt the widely-used spike-and-slab prior proposed in [51], the model prior p_δ will be generated using a weighted sum of two Gaussian distribution – a wide (slab) Gaussian and a narrow (spike) Gaussian:

$$q_\delta(\delta) = \prod_i q_{\delta_i}(\delta_i) = \prod_i \frac{q_{\text{slab}}(\delta_i) + \alpha q_{\text{spike}}(\delta_i)}{1 + \alpha} \quad (19)$$

with

$$q_{\text{slab}}(\delta_i) = \mathcal{N}(\delta_i|0, \sigma^2), \quad q_{\text{spike}}(\delta_i) = \mathcal{N}\left(\delta_i|0, \left(\frac{\Delta}{6}\right)^2\right), \quad (20)$$

where $\alpha \geq 0$ is a hyperparameter to determine the height of the spiky Gaussian with respect to the wider slab. Given q_δ , the model prior p_δ is derived as

$$\begin{aligned} p_\delta(\delta) &= \prod_i \underbrace{\left(q_{\delta_i}(\delta_i) * \left(\Delta \cdot \mathcal{U}\left(-\frac{\Delta}{2}, \frac{\Delta}{2}\right) \right) \right)}_{p_{\delta_i}}(\delta_i) \\ &= \prod_i \int_{\delta_i - \frac{\Delta}{2}}^{\delta_i + \frac{\Delta}{2}} q_{\delta_i}(\delta'_i) d\delta'_i. \end{aligned} \quad (21)$$

The discrete model prior $p_{\delta_i}(\bar{\delta}_i)$ is the pushforward of p_{δ_i} by substituting $\bar{\delta}_i$ into (21), which equals to the mass of GMM

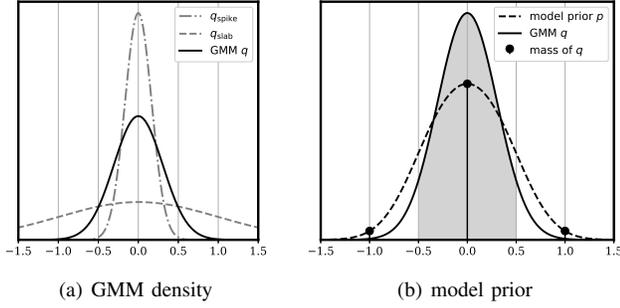


Fig. 5. A visual example of our GMM density function. The density q_{slab} is generated with $\sigma = 1$, and the quantization bin width is set to $\Delta = 1$ such that the standard deviation q_{spike} is $\frac{1}{6}$. The height of the spiky is $\alpha = 5$. Our model prior p_{δ_i} provides a continuous function, each value $p_{\delta_i}(\bar{\delta}_i)$ equals to the mass of GMM density $q_{\delta_i}(\bar{\delta}_i)$ in the quantization bin Δ centered at $\bar{\delta}_i$.

density $q_{\delta_i}(\bar{\delta}_i)$ in the quantization bin Δ centered at $\bar{\delta}_i$. By setting the standard deviation of the spike as $\frac{\Delta}{6}$, we ensure the mass within the zero-centering bin $[-\frac{\Delta}{2}, \frac{\Delta}{2}]$ cover 99.7% of the total mass (“ 3σ -criterion”). A visual illustration of GMM function and model prior is shown in Fig. 5.

The introduce of spike Gaussian distribution enforces sparsity on f_d and g_s model updates. A high spike weight (large α) results in almost negligible bandwidth cost for transmitting the zero-update model stream. By this means, the full-model adaptation procedure can learn to make a binary decision: a parameter is worth updating by using higher bandwidth cost, or it is not updated by transmitting the negligible zero-update sign. Note that the whole online learning process takes place in the transmitter since we assume the transmitter has a local copy of the decoder parameters, thus the above binary decision is also made inside the transmitter. When the entropy coding is applied on $\bar{\delta}_i$, zero-update corresponds to a sequence of negligible ‘spent’ bits. This trick can efficiently reduce the unnecessary model stream bandwidth cost.

The iterative transmitter model parameter updating procedure is the same as (13) by replacing $\mathcal{L}_{\text{R-D}}$ with $\mathcal{L}_{\text{R-D-M}}$, and the receiver model parameter updating procedures are written as

$$\theta_g^{(t)} = \theta_g^{(t-1)} - \gamma \nabla_{\theta_g} \mathcal{L}_{\text{R-D-M}}(\phi, \theta, \psi, \delta, \mathbf{x}, \mathbf{h}), \quad (22a)$$

$$\theta_f^{(t)} = \theta_f^{(t-1)} - \gamma \nabla_{\theta_f} \mathcal{L}_{\text{R-D-M}}(\phi, \theta, \psi, \delta, \mathbf{x}, \mathbf{h}). \quad (22b)$$

Since the model stream bandwidth cost M is much smaller than the data stream bandwidth cost R , the transmission of model updates δ default to a classical digital communication link, which will be quantized, entropy coded, channel coded, and modulated as a digital symbol sequence \mathbf{m} passing over the wireless channel W . This transmission process is assumed to be reliable. Thus, the received $\hat{\delta}$ is the quantized version $\bar{\delta}$, and the model stream rate should be computed by substituting the uniformly-noised proxy $\bar{\delta}$ into the model prior p_{δ} , i.e., $-\log p_{\delta}(\bar{\delta})$, in the R-D-M loss function $\mathcal{L}_{\text{R-D-M}}$. The encoding and decoding procedure of the transceiver full-model online learned NTSCC system is shown in Fig. 6. The whole online adaptation procedure is defined formally in Algorithm 3.

Algorithm 3: TxRx Full-Model Online Adaptation

Input: Baseline model parameters (ϕ, θ, ψ) trained on training set, model parameter quantizer Q_{Δ} , model prior p_{δ} , the source dataset to be overfitted \mathcal{X} , the CSI vector \mathbf{h} .

Output: The updated encoder models ϕ_g^* and ϕ_f^* and quantized decoder model updates $\bar{\delta}$.

```

1 procedure ONLINEUPDATEMODELS( $\phi, \theta, \psi, \mathbf{x}, \mathbf{h}$ )
2   Initialize model parameters:  $\phi_g^{(0)} \leftarrow \phi_g,$ 
    $\phi_f^{(0)} \leftarrow \phi_f, \theta_g^{(0)} \leftarrow \theta_g, \theta_f^{(0)} \leftarrow \theta_f;$ 
3   for  $t = 1, 2, \dots, T_{\max}$  do
4     Quantize transmittable model parameters:
        $(\bar{\theta}_g^{(t-1)}, \bar{\theta}_f^{(t-1)}) \leftarrow \bar{\delta} + (\theta_g, \theta_f)$ , with
        $\bar{\delta} = Q_{\Delta}(\delta)$  and
        $\delta = (\theta_g^{(t-1)}, \theta_f^{(t-1)}) - (\theta_g, \theta_f);$ 
5     Draw a batch of data samples  $\mathbf{x}$  from the
       source domain  $\mathcal{X}$ ;
       //specifically used for domain
       adaptation
6     Forward pass:  $\mathbf{x} \xrightarrow{g_a(\cdot; \phi_g^{(t-1)})} \mathbf{y} \xrightarrow{f_e(\cdot; \phi_f^{(t-1)})}$ 
        $\mathbf{s} \xrightarrow{W(\cdot; \mathbf{h})} \hat{\mathbf{s}} \xrightarrow{f_d(\cdot; \bar{\theta}_f^{(t-1)})} \hat{\mathbf{y}} \xrightarrow{g_s(\cdot; \bar{\theta}_g^{(t-1)})} \hat{\mathbf{x}}$  with
       the latent prior  $\mathbf{y} \xrightarrow{h_a(\cdot; \phi_h)} \mathbf{z} \xrightarrow{h_s(\cdot; \theta_h)} \{\boldsymbol{\mu}, \boldsymbol{\sigma}\};$ 
7     Compute loss  $\mathcal{L}_{\text{R-D-M}}$  according to (16) by
       setting  $\hat{\delta} = \bar{\delta}$  and the model rate proxy
        $-\log p_{\delta}(\bar{\delta})$ ;
8     Backpropagate using STE for  $Q_{\Delta}$ , and update
        $(\phi_g^{(t)}, \phi_f^{(t)}, \theta_g^{(t)}, \theta_f^{(t)})$  using gradients
        $(\nabla_{\phi_g} \mathcal{L}_{\text{R-D-M}}, \nabla_{\phi_f} \mathcal{L}_{\text{R-D-M}}, \nabla_{\theta_g} \mathcal{L}_{\text{R-D-M}}, \nabla_{\theta_f} \mathcal{L}_{\text{R-D-M}})$ 
       as (13) and (22);
9   return Updated models  $\phi_g^* \leftarrow \phi_g^{(T_{\max})}$  and
        $\phi_f^* \leftarrow \phi_f^{(T_{\max})}$ , and quantized model updates  $\bar{\delta}$ .
```

IV. OVERFITTING THE CHANNEL

In this section, we discuss how to adapt to the instant CSI $\mathbf{h} \in \mathcal{H}'$ such that the channel domain adapted model in the previous Section III can be further upgraded to the instance adapted model, i.e.,

$$\Sigma_{(\mathbf{x}, \mathcal{H}')} \Rightarrow \Sigma_{(\mathbf{x}, \mathbf{h})}, \quad (23a)$$

$$\Sigma_{(\mathcal{X}', \mathcal{H}')} \Rightarrow \Sigma_{(\mathcal{X}', \mathbf{h})}. \quad (23b)$$

It can be observed that the value of instant \mathbf{h} is a key factor affecting the performance of deep JSCC codec. Under different channel states, different resource allocation strategies should be adopted to implicitly adjust source coding rate and channel coding rate inside the JSCC codec. The vanilla NTSCC model [24] is trained under the objective (10) using different signal-to-noise ratio (SNR) indicating different channel states. Recently, an attention module was introduced into convolutional neural networks (CNN) to make deep JSCC SNR-adaptive using a single model [20]. However, a simple average SNR cannot fully capture all the characteristics of the channel state

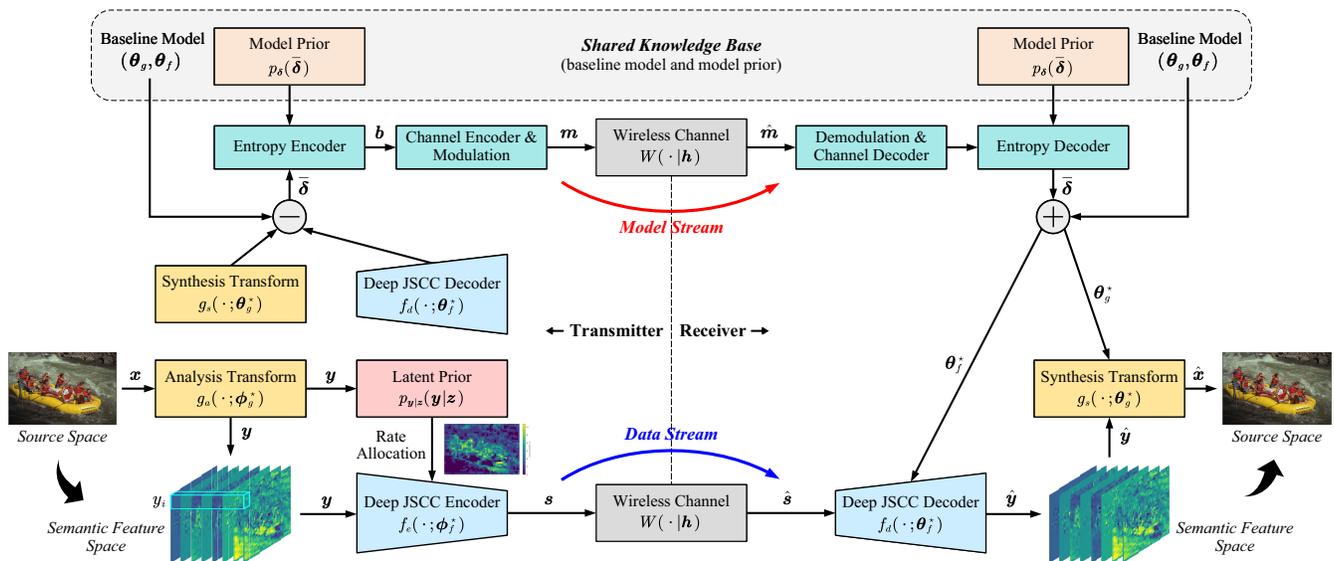


Fig. 6. Visualization of encoding and decoding of our full-model online learned NTSCC system.

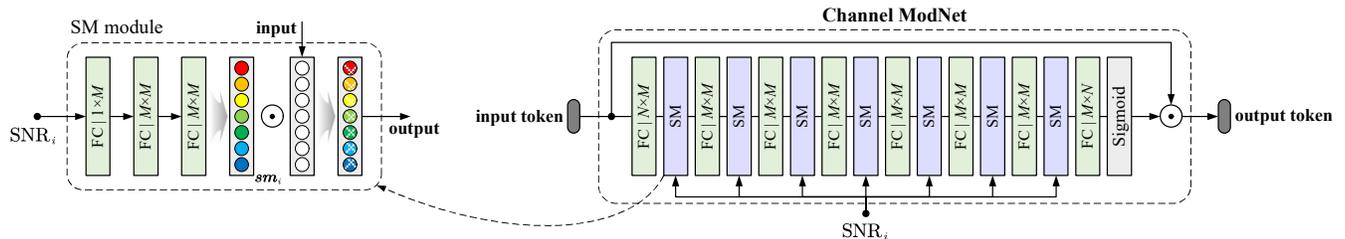


Fig. 7. Channel-dependent deep JSCC codec used for overfitting the CSI instance in our adaptive NTSCC system. FC denotes fully-connected network, and the following parameter marks the “input dimension \times output dimension”. “ \odot ” denotes the element-wise product.

\mathbf{h} . To tackle this, a CA-JSCC scheme was proposed in [52] that employs dual-attention mechanism to achieve adaptation over multipath OFDM channels. Even though these methods have proposed novel adaptation strategies, they are mostly based on traditional convolutional auto-encoders. In this paper, we investigate a new channel-dependent mechanism specifically tailored for vision Transformers (ViTs) adopted in NTSCC, enabling the whole system automatically to adapt to \mathbf{h} without relying on gradient descent.

Our idea is introducing a *plug-in* module to modulate the output of ViT-based deep JSCC encoder in NTSCC as shown in Fig. 7. The proposed “Channel ModNet” is a plug-in module inserted as the last layer of deep JSCC encoder or the first layer of deep JSCC decoder. For different channel states \mathbf{h} , we obtain different neural-syntax to generate a more specific deep JSCC codec functions f_e and f_d .

In particular, the architecture of the proposed Channel ModNet for instant channel state \mathbf{h} adaptation is depicted in Fig. 7. ModNet consists of 8 FC layers separated by 7 SNR modulation (SM) modules. SM is a three-layered FC network with input being the channel SNR that corresponds to y_i (denoted as SNR_i). As stated before, the receiver can obtain the instant CSI \mathbf{h} via channel estimation, thus the ModNet in the deep JSCC decoder f_d can obtain the explicit SNR_i for each y_i . Herein, SNR_i is computed by averaging the channel symbol SNRs along the \bar{k}_i -dimensional sequence s_i . To be

aligned with practical systems, we assume the CQI available at the transmitter via a feedback link, which represents an averaged SNR along all transmitted symbols s . Thus, the SM module input is $\text{SNR}_i = \overline{\text{SNR}}$ at the transmitter. The SM module transforms the input SNR_i into an c -dimensional tensor sm_i as depicted in Fig. 7. In this manner, the arbitrary target modulator can be realized by assigning a corresponding SNR value. And the CSI \mathbf{h} is therefore associated with the c -dimensional tensor sm_i in each SM module. Then, the input c -dimensional feature will be fused with sm_i in the element-wise product.

Multiple SM modules are cascaded sequentially in a coarse-to-fine manner, as shown in Fig. 7. The previously modulated features are fed into subsequent SM modules. In our design, the channel state modulation is comprehensively considered in a token-wise attention fashion. In this way, our channel-dependent modulation method can notify the JSCC codec what the current channel state is. The proposed channel-dependent method can be applied to arbitrary wireless channel model, and the ModNet trained under a specific channel domain \mathcal{H}' can provide better channel adaptation performance for the instances $\mathbf{h} \in \mathcal{H}'$.

V. EXPERIMENTAL RESULTS

In this section, we first outline the experimental configurations and evaluation protocol, followed by the presentation

of source/channel overfitting evaluation results for the online learned NTSCC. Also, we carry out a comprehensive ablation study and provide analytical discussions to highlight the key advantages and clarify some limitations of this work.

A. Experimental Setup

1) *Datasets*: Regarding the wireless video transmission, we adopt the Video Streaming Dataset 4K (VSD4K) dataset [36] for evaluation. The VSD4K dataset is composed of several 4K 30fps videos, including six popular video categories: game, vlog, interview, sport, dance, and city. The first three categories are mainly single-scene but have multiple points of view, while the latter contains various scenes and large-scale motion. In this paper, we consider the I-frame wireless transmission problem, i.e., the test set consists a group of I-frame images from a single video of a specific category, which is a subproblem of video transmission. We select a 45s-length representative video from each category, subsample each video by the first of four frames to create a dataset of I-frames, and generate low-resolution videos of 480P, 720P, and 1080P using bicubic interpolation. Unless otherwise specified, the following experiments will be carried out on 480P video sources.

2) *Network Implementation*: Note that the proposed source and channel overfitting method are model agnostic that can be transplanted to any neural network based end-to-end transmission system, we employ NTSCC [24] as a representative architecture for online updating in this paper. As stated before, we do not transmit the side information \bar{z} to simplify the implementation of NTSCC. We use the DIV2K [53] image set to train our baseline model, which contains 800 natural images of 2K resolutions on average. The baseline model training procedure includes 6,000 epochs using Adam optimizer [54] with the learning rate 1×10^{-4} . During model training, the images are randomly cropped to 256×256 patches to form a batch of 8 patches. In all experiments, we train four baseline models with $\lambda = 256, 64, 16, 4$ respectively under the additive white Gaussian noise (AWGN) channel, where the scaling factor η_y is set to 0.2 at SNR = 10dB and 0.4 at SNR = 0dB. The mean squared error (MSE) is used as the distortion measurement to qualify the end-to-end image transmission performance. For subjective comparison, we further consider perceptual quality optimization to better align with the human vision in semantic communication system. The training details of NTSCC (Perceptual) are the same as that in [24].

3) *Transmitter Adaptation Details*: The Tx-adaptation schemes comprise two main approaches, namely the online transmitter model updating and the transmitter code updating. The former can be applied either to every I-frame (i.e., instance adaptation) or to every video sequence (i.e., domain adaptation) within the VSD4K dataset. The latter is applicable only in the instance adaptive mode. Specifically, for the instance adaptation, the transmitter model adaptation scheme updates (g_a, f_e) for $T_{\max} = 100$ steps per I-frame with the learning rate $\gamma = 1 \times 10^{-4}$. Regarding transmitter code adaptation, we first update \mathbf{y} for $Y_{\max} = 50$ steps and then update \mathbf{s} for another $S_{\max} = 50$ steps, where the learning rate is set to

$\gamma = 1 \times 10^{-3}$ in both optimization procedures. It is possible to update with more steps using the same or lower learning rate, but our experiments have shown that negligible performance gain can be obtained.

4) *Transceiver Adaptation Details*: The TxRx-adaptation scheme builds upon the aforementioned pre-trained 4 baseline models, we online update the parameters $(\phi_g, \phi_f, \theta_g, \theta_f)$ for the R-D-M loss $\mathcal{L}_{\text{R-D-M}}$ over each video sequence. To ensure the average channel bandwidth cost of updated models close to the baseline, the R-D-M trade-off hyperparameter tuples (λ, β) is set to (256, 4), (32, 4), (4, 1), and (2, 1), which correspond to the baseline models learned with $\lambda = 256, 64, 16, 4$, respectively. For the decoder updating process, which will be transmitted to the receiver as the model stream, we set the parameter quantization bin width as $\Delta = 0.005$ and utilize the GMM prior, where σ in q_{slab} is set to 0.05, and the spike weight $\alpha = 1000$. The number of updating steps is set to $T_{\max} = 10000$, with the learning rate of $\gamma = 1 \times 10^{-5}$ for g_a and f_e , and $\gamma = 1 \times 10^{-4}$ for g_s and f_d . We select the best updated model based on the R-D-M loss for model inference.

5) *Comparison Schemes*: Apart from the comparison with baseline NTSCC system, we also compare our online learned NTSCC with SOTA engineered wireless video transmission systems. Specifically, we compare our approach with HEVC (H.265) [28] and VVC (H.266) [29] for video compression combined with 5G LDPC code [23] for channel coding. We use “+” operator to combine source coding and channel coding schemes for brevity. For VVC and HEVC, we use the official test model VTM-12.2 with intra profile and BPG software to encode I-frames in the YUV444 mode, and calculate the PSNR in RGB. Additionally, the NTC source compression scheme proposed in [43] combined with 5G LDPC codes are also included in comparison. We adopt the released NTC models based on mean and scale hyper-prior implemented by CompressAI [55]. To ensure a fair comparison, we crop all I-frames to multiples of 64 to avoid padding for neural codecs. The above simulations are implemented on the top of Sionna [56], which is a open-source library for link-level simulations of digital communication systems. Apart from these methods, we also compare our online learned NTSCC with the emerging neural deep JSCC transmission scheme [15]. The channel bandwidth cost counts up all the transmitting streams over the wireless channel, including both data and model in our full-model adaptation mode.

B. Objective Performance Analysis

Fig. 8(a)–(h) show the averaged R-D performance under AWGN channels with SNR = 0dB (poor channel condition) and SNR = 10dB (good channel condition), respectively. In Fig. 8, the legend “NTSCC (TxRx-adapt)” denotes NTSCC with our transceiver full-model domain adaptation, and “NTSCC (Tx-adapt)” denotes our transmitter model instance adaptation. For the classical separation-based schemes, according to the adaptive modulation coding (AMC) mechanism, we need to traverse given combinations of LDPC coded modulation schemes to identify the highest-efficiency scheme under the reliable transmission constraint (block error rate \leq

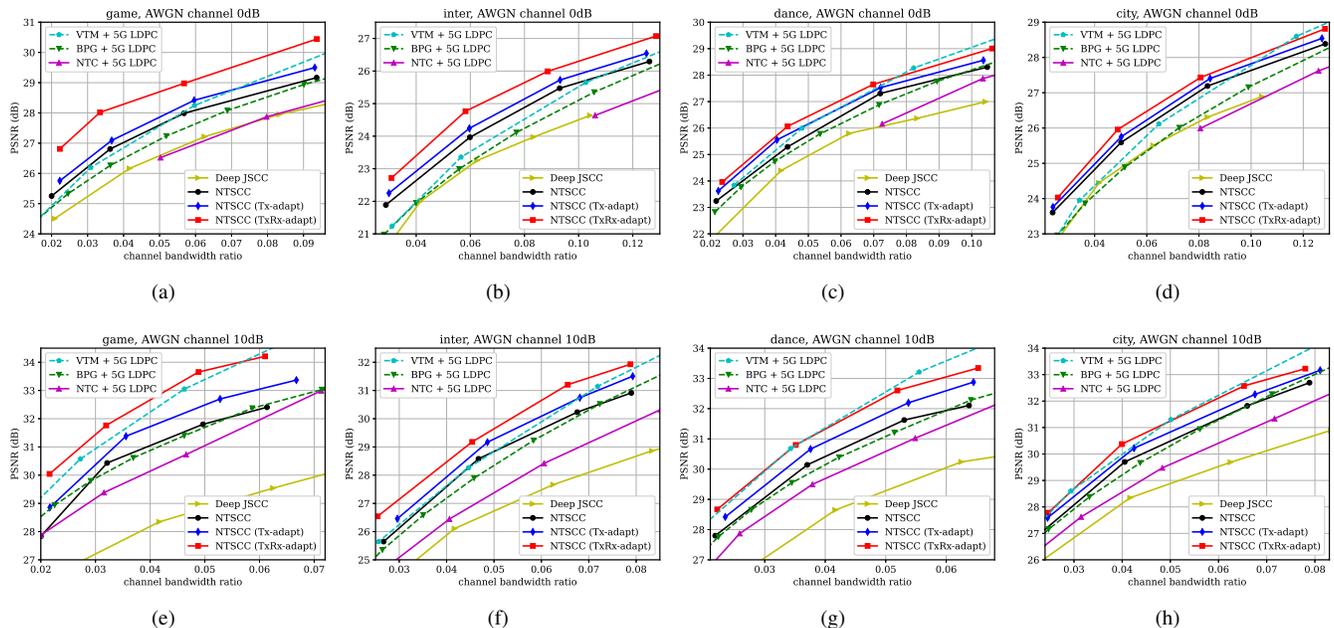


Fig. 8. PSNR performance versus the average channel bandwidth ratio (CBR) ρ over the AWGN channel at SNR = 0dB (bad channel condition, subfigures (a)–(d)) and SNR = 10dB (good channel condition, subfigures (e)–(h)). In this figure, our adaptive NTSCC systems are online updated toward each video sequence, which is sequence domain adaptation.

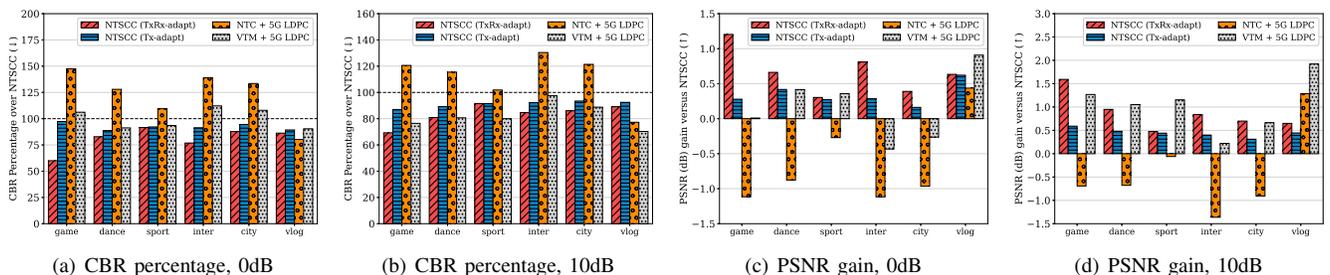


Fig. 9. Average CBR percentage and PSNR gains versus the baseline NTSCC semantic communication system over the AWGN channel at SNR = 0dB and SNR = 10dB, where “↓” on the Y-axis label denotes “lower is better”, and “↑” denotes “higher is better”.

10^{-5}). Accordingly, we adopt a $1/3$ rate (2048, 6144) LDPC code with 4QAM at SNR = 0dB, and a $2/3$ rate (4096, 6144) LDPC code with 16-ary quadrature amplitude modulation (16QAM) at SNR = 10dB. Results in Fig. 8 indicate that by using our proposed online overfitting mechanisms, we achieve considerable gains compared to the baseline NTSCC, and the transceiver full-model adaptation results in higher performance gain than updating the transmitter only. Moreover, compared to the best performed engineered transmission system “VTM + 5G LDPC”, our NTSCC (TxRx-adapt) shows comparable or even better performance. As one typical scheme towards adaptive semantic communications, our online learned NTSCC shows the potential becoming the emerging SOTA scheme for end-to-end wireless data transmission. The existing well-known improvements of semantic communication are compared with “BPG + 5G LDPC” [24] or on the tiny-scale image dataset [14], [15], our performance gain is more meaningful as the first end-to-end data transmission scheme overpassing the SOTA coded transmission system (VTM + 5G LDPC) on high-resolution image/video datasets.

We further plot the PSNR gain or bandwidth saving of the six video sequences in Fig. 9. Here, we employ the widely-used Bjøntegaard Delta (BD) rate reduction algorithm [57] for relative performance evaluation. The baseline scheme is NTSCC trained on the DIV2K dataset. Fig. 9(a) and 9(b) show the relative average CBR percentage over the baseline at the same PSNR under channel SNR 0dB and 10dB, respectively. Fig. 9(c) and 9(d) show the average quality improvement in terms of PSNR over the baseline for each scheme. Specifically, NTSCC with transceiver adaptation (TxRx-adapt) can save up to 41% bandwidth cost, while NTSCC with only transmitter adaptation (Tx-adapt) can save up to 13%. Among the six sequences in VSD4K, our NTSCC (TxRx-adapt) performs better than all other systems for most categories (4/6). However, it shows slightly worse performance on the *sport* and *vlog* sequences. We attribute this performance loss to the presence of multiple shots and diverse scenarios in these two sequences, which complicates the adaptation of the baseline NTSCC model to the entire video sequence. To improve performance, it is recommend to restrict the adaptation range to align with

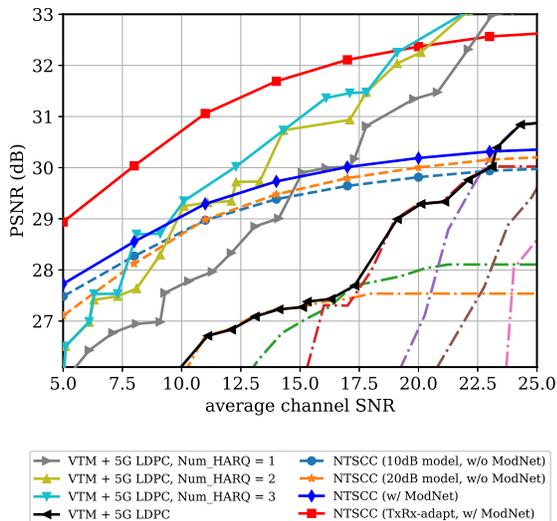


Fig. 10. PSNR performance versus the change of average SNR over the COST2100 5.3GHz indoor channel under the CBR constraint $\rho \leq 0.03$. The two baseline NTSCC models are trained under this COST2100 fading channel with average SNR = 10dB and average SNR = 20dB, respectively.

the obvious scene changes within a video sequence.

Next, we verify the flexibility of our proposed plug-in Channel ModNet on adaptation to the instant channel state. We consider the widely-used COST2100 wireless fading channel model [58] for evaluation. CSI samples are collected in an indoor scenario at 5.3GHz bands. In this case, the transmitted symbols s will pass over a frequency selective channel on the OFDM grid. We transmit the I-frames of *game* video sequence over this channel.

Fig. 10 shows the channel adaptation results. It can be seen that NTSCC with ModNet (see the solid blue curve) performs better than the standard NTSCC model (without ModNet) trained at a given average SNR. It verifies the Channel ModNet module can make NTSCC model aware of the instant token-wise SNRs, which is more accurate than using a rough average SNR among all embeddings.

Moreover, performance can be further improved by utilizing full-model adaptation based on the NTSCC with ModNet, as illustrated by the red line in Fig. 10. This approach enables both source domain adaptation and channel instance adaptation. Additionally, we also present a comparison with VTM + 5G LDPC schemes using the same CSI samples collected from COST2100. The configurations of LDPC codes and QAM are consistent with those in [15], and we take the envelope of all combinations of coded transmission schemes as the final performance (black line in Fig. 10). Apparently, it cannot provide satisfactory quality using one-shot transmission due to the cliff-effect. To align better with practical communication system, we further adopt the hybrid automatic repeat request (HARQ) with chase combining (CC) [59] to enhance the system performance while our NTSCC schemes continues to use one-shot transmission. It can be seen that retransmissions bring considerable PSNR gains especially for the high SNR region, and the gain increases with the number of HARQ allowed. However, HARQ also introduces much higher transmission

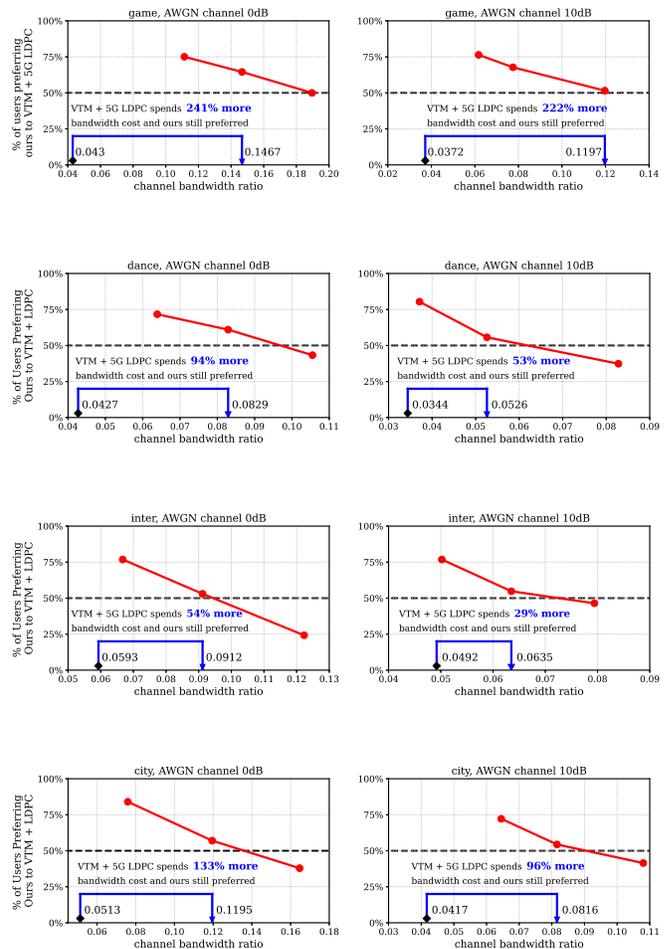


Fig. 11. User study results. As an anchor, the black diamond point denotes the CBR cost of our adaptive NTSCC (TxRx-adapt) semantic communication system optimized under perceptual loss. For each data point in the red line, its Y-axis shows what percentage of users prefer our scheme, and the X-axis is the CBR cost of VTM + 5G LDPC. The blue arrow highlights how much extra CBR cost VTM + 5G LDPC spends where still more than 50% of participants prefer ours. The results are counted from 2953 ratings in total, with an average of 123 per data point.

latency. As a comparison, the TxRx-adapt NTSCC (red line in Fig. 10) can provide competitive performance with one-shot transmission.

C. Visual Examples and User Study

Since the quality metrics sometimes fail to account for many nuances of human perception, we therefore rely on human opinions collected in a thorough user study. We employ the “two alternatives, forced choice” (2AFC) test for quantitative evaluation. As a widely-used approach in the perceptual visual quality assessment [6], [60], [61], the user study interface shows human raters a group of three images in the same column, where the middle image is always the original image, and the other two are lossily generated from two different methods. Then, we require raters to choose which image (left or right) looks closer to the original. We adopt the NTSCC (TxRx-adapt) model learned under perceptual loss [24] and

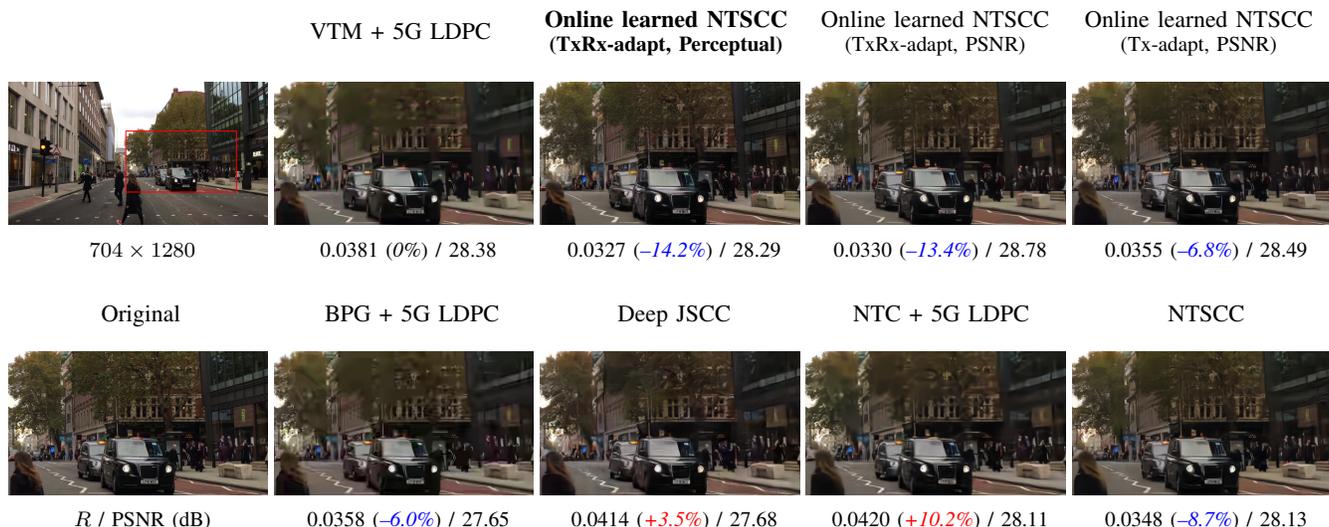


Fig. 12. Examples of visual comparison. The first column shows the original image and its cropped patch. The second to the fifth column show the reconstructed images by using different transmission schemes over the AWGN channel at SNR = 0dB, where the metrics in parentheses indicate the model training target loss function (PSNR or perceptual loss). Red number and blue number indicate the percentage of bandwidth cost increase and saving compared to the “VTM + 5G LDPC” scheme.

compare it to VTM + 5G LDPC with close or more CBR cost. In particular, we randomly collect 10 I-frame images for each video sequence transmitted over the AWGN channel at SNR = 0dB and SNR = 10dB. The entire rating process consists of several rounds in which all participants were asked to vote on each of the 10 groups of images. If our NTSCC (TxRx-adapt) dominates, the CBR cost of VTM + 5G LDPC will be increased in the next round until it receives more votes. User study results are given in Fig. 11. As we can see, our source and channel overfitting approach enables NTSCC to greatly outperform the SOTA VTM + 5G LDPC transmission scheme, especially for `game` and `city` sequences. It verifies that the proposed online learned NTSCC system can better support the human perceptual vision demands in wireless data transmission, which is aligned with the target of semantic communications.

To intuitively demonstrate the effect of the proposed overfitting method, we further pick visible results on transmitting the `city` video sequence in Fig. 12. From these results, we can observe that online learned NTSCC optimized under the perceptual loss achieves much better visual quality than other schemes with lower channel bandwidth cost.

D. Ablation Study

1) *R-D-M Tradeoff*: Table I provides an ablation study concerning the impact of β in domain adaptation, which investigates the tripartite trade-off among the data stream bandwidth cost (R), model stream bandwidth cost (M), and end-to-end distortion (D). Specifically, we demonstrate the distribution of total CBR cost using a fixed $\lambda = 4$ in a 45s `game` sequence (including $N = 372$ I-frames) over the AWGN channel at SNR = 10dB. Herein, the model stream cost M is counted according to the channel state and the AMC mechanism, as that used for the traditional separation-based schemes in Fig. 8. When the LDPC coding rate and digital

TABLE I
R-D-M TRADE-OFF AS A FUNCTION OF β , WHERE THE TOTAL CBR $\rho = (R + M/N)/m$.

| β | PSNR (dB) | ρ (total CBR) | R/m (data CBR) | $M/(mN)$ (model CBR) |
|---------|-----------|--------------------|------------------|----------------------|
| – | 31.80 | 0.04967 | 0.04967 | 0 |
| 16 | 33.13 | 0.04870 | 0.04860 | 0.00011 |
| 4 | 33.33 | 0.04853 | 0.04825 | 0.00027 |
| 1 | 33.65 | 0.04885 | 0.04811 | 0.00073 |
| 0.1 | 34.35 | 0.05245 | 0.04867 | 0.00378 |
| 0.01 | 35.21 | 0.07377 | 0.04549 | 0.02828 |

modulation order are determined, the number of transmitted symbols M can be calculated according to the estimated model stream entropy in (16), (21) and the selected AMC scheme. From this table, we observe that as the value of β decreases, the R-D terms start to dominate the R-D-M trade-off, resulting in improved R-D performance compared to the baseline model (in the first row of Table I) and an increased bandwidth cost for the model stream. since the decoder is updated only once to adapt to the `game` domain, the bandwidth cost of the model stream can be amortized over a substantial number of I-frames. As a result, the model rate becomes negligible, especially in the high β region. An appropriate selection of β in the transceiver adaptation can contribute to a significant improvement in reconstruction quality while also reducing the total bandwidth cost $k = R + \frac{M}{N}$.

We extend our analysis by examining the impact of λ in domain adaptation. Fig. 13 depicts the model updating progress on the CBR-PSNR plot with various λ values and a fixed $\beta = 1$. From the results, we observe that the significant gains from domain adaptation are achieved during the early stages of model updating, within just hundreds of iterations,

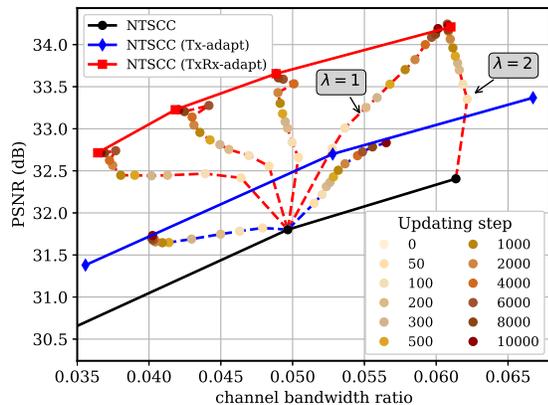


Fig. 13. Online learning on the *game* video under AWGN channel at SNR = 10dB. Each dot shows performance under different updating steps.

and the performance continues to improve moderately during online updating. Additionally, optimizing for different λ values can cause the R-D-M curve to move in different directions, but the final R-D-M performance point seems independent of the starting point from baseline model, as seen in the highlighted $\lambda = 1$ and $\lambda = 2$ curves. This phenomenon suggests that, given the model stream weight β , the maximum reduction in the amortization gap on the receiver has almost been reached. To further reduce the amortization gap, the only approach is to limit the range of the adaptation domain.

2) *Performance-Complexity Tradeoff*: We present the average end-to-end processing latency and the coding efficiency of considered coded transmission schemes on the 480p VSD4K dataset in Table II. This includes coding latency and the BD- ρ (“ \downarrow ” means smaller is better), which measures bandwidth ratio savings at the equivalent quality, “+/-” signs indicate more/fewer bandwidth cost than baseline [57]. All experiments are conducted using PyTorch 1.9.0, with the Inter Xeon Gold 6226R CPU (mainly for arithmetic codec) and one RTX 3090 GPU (for model updating and inference, and LDPC codec).

Results show that learned end-to-end data transmission schemes (e.g., Deep JSCC and NTSCC) run much faster than classical layered designed schemes. This makes emerging end-to-end data transmission schemes particularly suitable for low-latency applications such as live video delivery, which requires real-time encoding on the client side. We then discuss the online learning schemes. Essentially, they enhance the coding efficiency of end-to-end data transmission schemes by leveraging computation at the transmitter. For the instance adapt scheme, it requires tens of seconds of additional encoding time per frame for updating codes or encoder model, in exchange for 9% bandwidth saving compared to the standard NTSCC. The instance adapt paradigm can potentially support wireless image/video delivery services, where encoding can be run in parallel and in advance.

For the TxRx domain adaptation scheme, in order to achieve improved coding efficiency, it is necessary to perform online updating of the full-model parameters and transmit the model stream to the receiver before initiating the encoding process. Consequently, there will be a delay before the full-model

TABLE II
TRANSMISSION EFFICIENCY AND AVERAGED ENCODING/DECODING LATENCY COMPARISON.

| Transmission scheme | BD- ρ (% , \downarrow) | | End-to-end latency | |
|---------------------|--------------------------------|-------|--------------------|-------|
| | 0dB | 10dB | Enc. | Dec. |
| BPG + 5G LDPC | 13.1 | 1.1 | 3.8s | 380ms |
| VTM + 5G LDPC | -1.2 | -17.7 | 114s | 430ms |
| NTC [43] + 5G LDPC | 22.9 | 11.2 | 27ms | 60ms |
| Deep JSCC [15] | 42.3 | 46.4 | 17ms | 25ms |
| NTSCC [24] | 0 | 0 | $t_0 = 25$ ms | 15ms |
| Tx-adapt, instance | -9.0 | -9.0 | t_1 | 15ms |
| TxRx-adapt, domain | -19.2 | -16.4 | t_2 or t_0 | 15ms |

- ¹ The Tx-adapt encoding time t_1 is associated with the total number of instance adaptation steps. As a reference, it takes about 29.7s for Tx model updating in Algorithm 1 with $T_{\max} = 100$ steps and 18.3s for updating latents in Algorithm 2 with $Y_{\max} = S_{\max} = 50$ steps.
- ² The TxRx-adapt online learned NTSCC system requires online updating full-model parameters before inference. As a reference, it takes about $t_2 = 2.2$ h for $T_{\max} = 10000$ steps. After updating, the subsequent encoding latency is same as the standard NTSCC, i.e., $t_0 = 25$ ms.

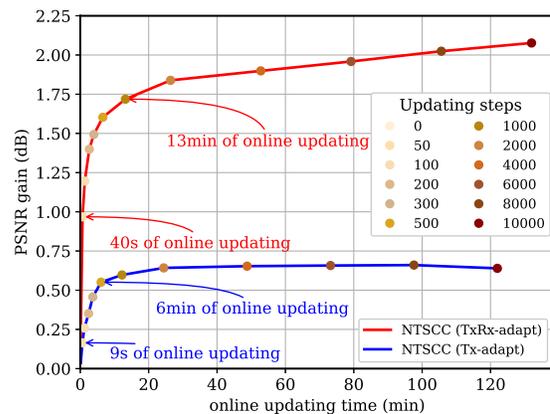


Fig. 14. The PSNR gain versus standard NTSCC over time during model updating. Each dot shows performance under different updating steps.

updating is completed. If the updating process is unacceptable for delay-sensitive applications (e.g., live video delivery), we recommend reverting to the baseline model for a while before the updating process is completed. In this way, one can ensure low end-to-end latency and achieve additional gains from domain adaptation in the same time.

Another advantage of the domain adaptation paradigm is that online fine-tuned model can converge more rapidly if the test-time content is relatively consistent, such as the scenarios including a fixed-location camera, cityscapes, live broadcasts, or video conferences. In Fig. 14, we plot the detailed PSNR gain versus the baseline over the online training time on the *game* sequence. As seen in this figure, about 50% of the gain can be achieved after updating for less than 1 minute, and 80% of the gain can be attained with about 10 minutes of updating. Since this is a “one-time” cost amortized across all the video frames, the additional training cost is negligible compared to the channel bandwidth saving benefits it brings, especially for long videos.

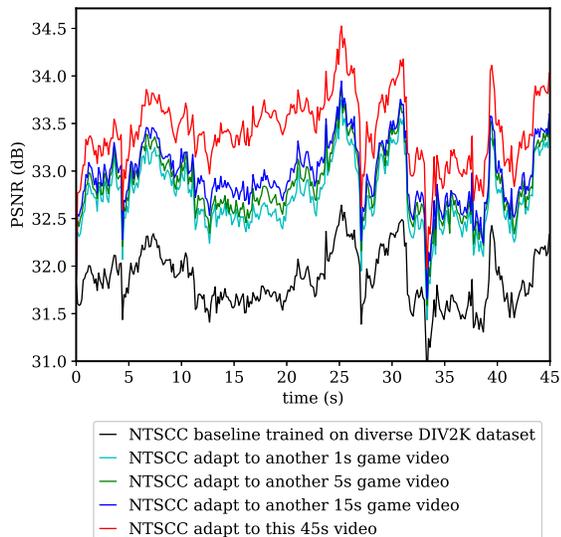


Fig. 15. Case study: live video delivery. We downloaded a new 45s video of the same game from YouTube and used it to simulate a live video delivery scenario, where the video content to be encoded cannot be known in advance. The bandwidth costs of different schemes are adjusted to be the same.

3) *Results for Unseen Data:* To evaluate the performance of our full-model adaptation method on unseen data, we downloaded a new 45s video from YouTube as test sequence, which contains the same game as the `game` sequence in the VSD4K dataset. We use the new video as unseen data to simulate a live video delivery scenario, where the video content to be encoded cannot be known in advance. In other words, the online learning can only utilize previously collected video content. The results in Fig. 15 demonstrate that the proposed online learning method can yield significant performance gains compared to the baseline, even when the training set consists of a short video sequence with a duration of only 1s. Moreover, the models adapt to another game video can still approximate the performance of the best models adapted to this precise video. This observation indicates that the proposed method exhibits good robustness to previously unseen data from the same category, as the model assimilates domain-specific knowledge through the online learning process.

4) *Result for Various Video Resolution:* We further report the performance on the `game` video of different resolutions in Fig. 16. These results validate the stability of quality improvements achieved using our proposed overfitting mechanism. The proposed online learned NTSCC shows superior performance over the SOTA VTM + 5G LDPC scheme, particularly for high-resolution video I-frames.

E. Results Discussion

Regarding the above experimental results and ablation study, we can draw two conclusions:

- First, our online overfitting mechanism enables semantic communication system to efficiently adapt to both source content and wireless channel state. It can greatly reduce the model amortization gap, which finally contributes to the system performance gain.

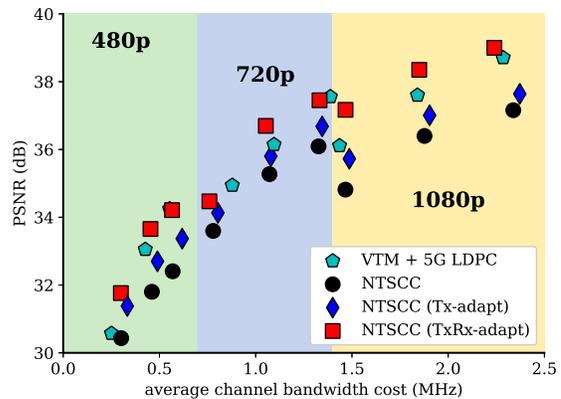


Fig. 16. PSNR performance of transmitting the `game` video at difference resolutions over the AWGN channel at SNR = 10dB, where the average bandwidth cost is computed within 100ms according to the 5G OFDM configurations: 14 OFDM symbols@1ms with the subcarrier bandwidth 15KHz.

- Second, the small additional complexity caused by model adaptation involves only at the transmitter. There is no extra complexity at the receiver, which is friendly to most content delivery tasks with the request of low decoding latency. It is aligned with the evolution idea of both traditional and neural video codec.

In summary, as mentioned in Section I, almost all traditional source compressors follow the hybrid transform coding paradigm to evolve, e.g., mode selection in HEVC [28] and VVC [29]. Accordingly, the idea of signal-dependent transform in traditional source compression methods inspires us to upgrade the deep learning based semantic communication system to a content-channel-dependent adaptive mode. The marriage of deep learning method and traditional coding idea can bring significant performance improvement.

Based on previous works about semantic communications [14]–[20], [24], [62], [63], and the work in this paper, we think that the advantages of neural network based end-to-end semantic communication system are three folds:

- First, the excellent content and semantics adaptivity of neural network is superior to signal processing based traditional models since the network parameters are learned based on lots of practical source and channel data samples while the models in the SOTA coded transmission standards are handcrafted based on prior knowledge.
- Second, the neural network can well represent and utilize source and channel features, which makes the semantic communication system can be optimized toward both human view perception and machine vision tasks. However, the existing source and channel coding standards only pursue high performance toward the objective quality assessment indices, e.g., PSNR.
- Third, the R-D optimization guided neural network training and adaptive switching for semantic communication is quite effective and efficient. As analyzed in this paper, a single model to deal with all the source data with diverse structures and varying wireless channel states is inefficient obviously. Therefore, the adaptively learning and switching according to source data and channel state is a

necessary solution to enhance the R-D performance for all deep learning based semantic communication systems. In addition, compared with the adaptive coding paradigm used in traditional source and channel coding where one appropriate coding mode is selected from the predefined mode pool with limited number of options, this paper leverages the overfitting property of neural network to adapt to arbitrary coding mode. Hence, the proposed adaptive semantic communication system is much more flexible to be tailored for specific source and channel instance with lower complexity.

In a nutshell, our overfitting method has efficiently catalyzed semantic communication system to provide more promising results. It is a necessary approach for all learning-based end-to-end communication system to further boost performance, which is indeed aligned with the evolution route of both traditional and neural video codec [64].

VI. CONCLUSION

In this paper, we have proposed an online learned end-to-end data transmission paradigm by well leveraging the deep learning model's overfitting property. Our adaptive NTSCC system can for instance or domain be updated after deployment, which leads to substantial gains on the bandwidth ratio-distortion performance. In this way, the emerging semantic communication is further upgraded to the adaptive semantic communication system. Specifically, we have proposed several lightweight methods to update the learned data transmission models. The ingredients of wireless transmitted stream include both the semantic representations of source data and the updated decoder model parameters. Accordingly, we have formulated the new optimization problem whose goal is minimizing the loss function that is a tripartite trade-off among data stream rate, model stream rate, and end-to-end distortion terms. Results have verified the substantial gains of the online learned NTSCC system, whose performance has surpassed the SOTA engineered coded transmission systems. As a new paradigm, our model adaptation methods have the potential to catalyze semantic communication upgrading to a new era.

REFERENCES

- [1] D. Gündüz, Z. Qin, I. E. Aguerri, H. S. Dhillon, Z. Yang, A. Yener, K. K. Wong, and C.-B. Chae, "Beyond transmitting bits: Context, semantics, and task-oriented communications," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 5–41, 2023.
- [2] P. Zhang, W. Xu, H. Gao, K. Niu, and et al., "Toward wisdom-evolutionary and primitive-concise 6G: A new paradigm of semantic communication networks," *Engineering*, vol. 8, pp. 60–73, 2022.
- [3] H. Xie, Z. Qin, G. Y. Li, and B.-W. Juang, "Deep learning enabled semantic communication systems," *IEEE Transactions on Signal Processing*, vol. 69, pp. 2663–2675, 2021.
- [4] J. Dai, P. Zhang, K. Niu, S. Wang, Z. Si, and X. Qin, "Communication beyond transmitting bits: Semantics-guided source and channel coding," *IEEE Wireless Communications, early access*, 2022.
- [5] C. E. Shannon, "A mathematical theory of communication, 1948," *Bell System Technical Journal*, vol. 27, no. 3, pp. 3–55, 1948.
- [6] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 586–595.
- [7] K. Ding, K. Ma, S. Wang, and E. P. Simoncelli, "Image quality assessment: Unifying structure and texture similarity," *arXiv preprint arXiv:2004.07728*, 2020.
- [8] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [9] L. Duan, J. Liu, W. Yang, T. Huang, and W. Gao, "Video coding for machines: A paradigm of collaborative compression and intelligent analytics," *IEEE Transactions on Image Processing*, vol. 29, pp. 8680–8695, 2020.
- [10] M. Fresia, F. Perez-Cruz, H. V. Poor, and S. Verdu, "Joint source and channel coding," *IEEE Signal Processing Magazine*, vol. 27, no. 6, pp. 104–113, 2010.
- [11] A. Guyader, E. Fabre, C. Guillemot, and M. Robert, "Joint source-channel turbo decoding of entropy-coded sources," *IEEE Journal on Selected Areas in Communications*, vol. 19, no. 9, pp. 1680–1696, 2001.
- [12] N. Ramzan, S. Wan, and E. Izquierdo, "Joint source-channel coding for wavelet-based scalable video transmission using an adaptive turbo code," *EURASIP Journal on Image and Video Processing*, vol. 2007, pp. 1–12, 2007.
- [13] C. Chen, L. Wang, and F. CM Lau, "Joint optimization of protograph LDPC code pair for joint source and channel coding," *IEEE Transactions on Communications*, vol. 66, no. 8, pp. 3255–3267, 2018.
- [14] E. Boursoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, 2019.
- [15] D. B. Kurka and D. Gündüz, "DeepJSCC-f: Deep joint source-channel coding of images with feedback," *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, pp. 178–193, 2020.
- [16] D. B. Kurka and D. Gündüz, "Bandwidth-agile image transmission with deep joint source-channel coding," *IEEE Transactions on Wireless Communications*, vol. 20, no. 12, pp. 8081–8095, 2021.
- [17] M. Jankowski, D. Gündüz, and K. Mikołajczyk, "Wireless image retrieval at the edge," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 89–100, 2020.
- [18] N. Farsad, M. Rao, and A. Goldsmith, "Deep learning for joint source-channel coding of text," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2018, pp. 2326–2330.
- [19] K. Choi, K. Tatwawadi, A. Grover, T. Weissman, and S. Ermon, "Neural joint source-channel coding," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2019, pp. 1182–1192.
- [20] J. Xu, B. Ai, W. Chen, A. Yang, P. Sun, and M. Rodrigues, "Wireless image transmission using deep source channel coding with attention modules," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 4, pp. 2315–2328, 2021.
- [21] T.-Y. Tung and D. Gündüz, "Deepwive: Deep-learning-aided wireless video transmission," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2570–2583, 2022.
- [22] Fabrice Bellard, "BPG image format," URL: <https://bellard.org/bpg/>.
- [23] Tom Richardson and Shrinivas Kudekar, "Design of low-density parity check codes for 5g new radio," *IEEE Communications Magazine*, vol. 56, no. 3, pp. 28–34, 2018.
- [24] J. Dai, S. Wang, K. Tan, Z. Si, X. Qin, K. Niu, and P. Zhang, "Nonlinear transform source-channel coding for semantic communications," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 8, pp. 2300–2316, 2022.
- [25] Z. Gu, W. Lin, B. Lee, and C. Lau, "Low-complexity video coding based on two-dimensional singular value decomposition," *IEEE Transactions on Image Processing*, vol. 21, no. 2, pp. 674–687, 2011.
- [26] C. Lan, J. Xu, W. Zeng, G. Shi, and F. Wu, "Variable block-sized signal-dependent transform for video coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 8, pp. 1920–1933, 2017.
- [27] S. Puri, S. Lasserre, and P. Le Callet, "Annealed learning based block transforms for HEVC video coding," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 1135–1139.
- [28] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (HEVC) standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [29] B. Bross, J. Chen, J.-R. Ohm, G. J. Sullivan, and Y.-K. Wang, "Developments in international video coding standardization after avc, with an overview of versatile video coding (VVC)," *Proceedings of the IEEE*, vol. 109, no. 9, pp. 1463–1493, 2021.
- [30] C. Cremer, X. Li, and D. Duvenaud, "Inference suboptimality in variational autoencoders," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1078–1086.

- [31] G. Lu, C. Cai, X. Zhang, L. Chen, W. Ouyang, D. Xu, and Z. Gao, "Content adaptive and error propagation aware deep video compression," in *European Conference on Computer Vision*. Springer, 2020, pp. 456–472.
- [32] J. Campos, S. Meierhans, A. Djelouah, and C. Schroers, "Content adaptive optimization for neural image compression," *arXiv preprint arXiv:1906.01223*, 2019.
- [33] Y. Yang, R. Bamler, and S. Mandt, "Improving inference for neural image compression," *Advances in Neural Information Processing Systems*, vol. 33, pp. 573–584, 2020.
- [34] T. Guo, J. Wang, Z. Cui, Y. Feng, Y. Ge, and B. Bai, "Variable rate image compression with content adaptive optimization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 122–123.
- [35] J. Kim, Y. Jung, H. Yeo, J. Ye, and D. Han, "Neural-enhanced live streaming: Improving live video ingest via online learning," in *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, 2020, pp. 107–125.
- [36] J. Liu, M. Lu, K. Chen, X. Li, S. Wang, Z. Wang, E. Wu, Y. Chen, C. Zhang, and M. Wu, "Overfitting the data: Compact neural video delivery via content-aware feature modulation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4631–4640.
- [37] T. Van Rozendaal, I. A.M. Huijben, and T. S. Cohen, "Overfitting for fun and profit: Instance-adaptive data compression," *arXiv preprint arXiv:2101.08687*, 2021.
- [38] D. Wang, W. Yang, Y. Hu, and J. Liu, "Neural data-dependent transform for learned image compression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17379–17388.
- [39] D. P Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [40] S. Wang, J. Dai, Z. Liang, K. Niu, Z. Si, C. Dong, X. Qin, and P. Zhang, "Wireless deep video semantic transmission," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 214–229, 2023.
- [41] J. Ballé, P. A Chou, D. Minnen, S. Singh, N. Johnston, E. Agustsson, S. J. Hwang, and G. Toderici, "Nonlinear transform coding," *IEEE Journal of Selected Topics in Signal Processing*, vol. 15, no. 2, pp. 339–353, 2020.
- [42] J. Ballé, V. Laparra, and E. P. Simoncelli, "End-to-end optimized image compression," in *Proceedings of the International Conference on Learning Representations*, 2016.
- [43] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," in *Proceedings of the International Conference on Learning Representations*, 2018.
- [44] I. H Witten, R. M Neal, and J. G Cleary, "Arithmetic coding for data compression," *Communications of the ACM*, vol. 30, no. 6, pp. 520–540, 1987.
- [45] J. Rissanen and G. Langdon, "Universal modeling and coding," *IEEE Transactions on Information Theory*, vol. 27, no. 1, pp. 12–23, 1981.
- [46] J. Ballé, N. Johnston, and D. Minnen, "Integer networks for data compression with latent-variable models," in *International Conference on Learning Representations*, 2018.
- [47] Y. Han, G. Huang, S. Song, L. Yang, H. Wang, and Y. Wang, "Dynamic neural networks: A survey," *arXiv preprint arXiv:2102.04906*, 2021.
- [48] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N Gomez, E. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [49] J. Jiang, V. Sekar, and H. Zhang, "Improving fairness, efficiency, and stability in http-based adaptive video streaming with festive," *IEEE/ACM Transactions on Networking*, vol. 22, no. 1, pp. 326–340, 2014.
- [50] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *arXiv preprint arXiv:1308.3432*, 2013.
- [51] V. Ročková and E. I George, "The spike-and-slab lasso," *Journal of the American Statistical Association*, vol. 113, no. 521, pp. 431–444, 2018.
- [52] H. Wu, Y. Shao, K. Mikołajczyk, and D. Gündüz, "Channel-adaptive wireless image transmission with ofdm," *IEEE Wireless Communications Letters*, vol. 11, no. 11, pp. 2400–2404, 2022.
- [53] E. Agustsson and R. Timofte, "Ntire 2017 challenge on single image super-resolution: Dataset and study," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 126–135.
- [54] D. P Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [55] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja, "CompressAI: a PyTorch library and evaluation platform for end-to-end compression research," *arXiv preprint arXiv:2011.03029*, 2020.
- [56] J. Hoydis, S. Cammerer, F. A. Aoudia, A. Vem, N. Binder, G. Marcus, and A. Keller, "Sionna: An open-source library for next-generation physical layer research," *arXiv preprint arXiv:2203.11854*, 2022.
- [57] G. Bjontegaard, "Calculation of average PSNR differences between RD-curves," *VCEG-M33*, 2001.
- [58] L. Liu, C. Oestges, J. Poutanen, K. Haneda, P. Vainikainen, F. Quitin, F. Tufvesson, and P. De Doncker, "The cost 2100 mimo channel model," *IEEE Wireless Communications*, vol. 19, no. 6, pp. 92–99, 2012.
- [59] A. Ahmed, A. Al-Dweik, Y. Iraqi, H. Mukhtar, M. Naeem, and E. Hossain, "Hybrid automatic repeat request (HARQ) in wireless communications systems and standards: A contemporary survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 4, pp. 2711–2752, 2021.
- [60] F. Mentzer, G. D Toderici, M. Tschannen, and E. Agustsson, "High-fidelity generative image compression," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [61] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. V. Gool, "Generative adversarial networks for extreme learned image compression," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 221–231.
- [62] S. Wang, J. Dai, S. Yao, K. Niu, and P. Zhang, "A novel deep learning architecture for wireless image transmission," in *2021 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2021, pp. 1–6.
- [63] M. Yang and H.-S. Kim, "Deep joint source-channel coding for wireless image transmission with adaptive rate control," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 5193–5197.
- [64] S. Ma, X. Zhang, C. Jia, Z. Zhao, S. Wang, and S. Wang, "Image and video compression with neural networks: A review," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 6, pp. 1683–1698, 2019.