# Towards Semantic Communications:
# Deep Learning-Based Image Semantic Coding

Danlan Huang, *Member, IEEE,* Feifei Gao, *Fellow, IEEE,* Xiaoming Tao, *Member, IEEE,* Qiyuan Du, and Jianhua Lu, *Fellow, IEEE*

## Abstract

Semantic communications has received growing interest since it can remarkably reduce the amount of data to be transmitted without missing critical information. Most existing works explore the semantic encoding and transmission for text and apply techniques in Natural Language Processing (NLP) to interpret the meaning of the text. In this paper, we conceive the semantic communications for image data that is much more richer in semantics and bandwidth sensitive. We propose an reinforcement learning based adaptive semantic coding (RL-ASC) approach that encodes images beyond pixel level. Firstly, we define the *semantic concept* of image data that includes the category, spatial arrangement, and visual feature as the representation unit, and propose a convolutional semantic encoder to extract semantic concepts. Secondly, we propose the image reconstruction criterion that evolves from the traditional *pixel similarity* to *semantic similarity* and *perceptual performance*. Thirdly, we design a novel RL-based semantic bit allocation model, whose reward is the increase in *rate-semantic-perceptual* performance after encoding a certain semantic concept with adaptive quantization level. Thus, the task-related information is preserved and reconstructed properly while less important data is discarded. Finally, we propose the Generative Adversarial Nets (GANs) based semantic decoder that fuses both locally and globally features via an attention module. Experimental results demonstrate that the proposed RL-ASC is noise robust and could reconstruct visually pleasant and semantic consistent image, and saves times of bit cost compared to standard codecs and other deep learning-based image codecs.

## Index Terms

Semantic communications, image semantic coding, Generative Adversarial Nets (GANs), reinforcement learning (RL), rate-semantic-perceptual criterion

D. Huang, F. Gao, X. Tao, Q. Du and J. Lu are with Beijing National Research Center for Information Science and Technology, Tsinghua University, Beijing, China. X. Tao is the corresponding author. Email:{huangdl, feifeigao, taoxm, lhh-dee}@mail.tsinghua.edu.cn.

## I. INTRODUCTION

Dating back to 1949, Weaver [40] defined three levels of communications: the bit level, the semantic level, and the effective level. The conventional communications falls into the bit level, where the receiver recovers the raw data by maximizing the accuracy of symbol transmission; The semantic communications interprets information at the semantic level and attempts to transmit the symbol that precisely conveys the desired meaning, instead of bit copy; The effective level concerns how effectively would the received meaning affect the behavior of the receiver in the desired way. Recently, deep learning achieves a breakthrough in semantic analysis task such as NLP, image processing, and speech recognition, etc, and therefore it paves the way for building the semantic communication systems and enables intelligent communication for human-to-machine as well as machine-to-machine.

Recent works mainly focus on semantic communications for text/speech modalities [35], [48]–[50], and therefore can be categorize as *text semantic communications*(TSC). A deep learning based joint source-channel coding (JSCC) of text is proposed in [14], which achieves lower word error rate and preserves semantic information of sentences. Later, the authors of [30] design an reinforcement learning (RL) based semantic communication system to deal with the non-differentiability of semantic metrics and interact with the surrounding noisy environment.

Actually, the vision modal is more informative than text modal, and the video/image traffic accounts for 75% of IP traffic nowadays [11]. However, the huge amount of multimedia traffic encounters the transmission challenges such as delay and network congestion, which are difficult to be solved by the conventional communication techniques. Hence, it is more important to build the *image semantic communications* (ISC) that can remarkably reduce the amount of data to be transmitted without sacrificing the semantic fidelity of the image.

The earlier attempt of semantic image communications [8], [25], [51] developed the JSCC method and achieved good performance in the challenging low signal-to-noise (SNR) and small bandwidth regimes. The work [8] maps the image pixel values directly to the complex-valued channel input symbols and learns noise resilient coded representations, and therefore outperforms separation-based digital communication at all SNR. Later, Kurka et.al [25] proposed an multiple-description JSCC scheme for bandwidth-agile image transmission. The attention DL based JSCC

for image transmission [51] successfully operates with different SNR levels during transmission.

However, the aforementioned works lack the interpretation ability of the image content. Actually, the main issue of the ISC is to discard the goal of precise reconstruction and pursue semantic fidelity of the reconstructed image even in aggressive compression ratio. The bit level image communications such as the standard codecs [1], [36], [47] and deep learning based codecs [23], [31], [33], [44], [45] aim to solve the rate-distortion optimization (RDO) problem. They exactly recover the transmitted image data at the receiver side by processing the image in pixel level, and therefore gradually approach the compression limits of Shannon's information theory. The highly compressed image essentially deviate from human perception and suffers from degradation such as blocking, ringing, blurry, and checkerboard artifacts [2], [22] that cause poor performance in semantic analysis tasks such as classification, detection, and segmentation [13], [17]. Therefore, blindly minimizing pixel-wise distortion may bring unnecessary bit overhead. The work [27] states that the conventional distortion metric Mean Squared Error (MSE) in compression cannot fulfill the requirements of desirable intelligent task performance, and it is essential to exploit machine-centric evaluation metrics for high inference accuracy.

Some preliminary studies have tried to leverage the semantic similarity as the reconstruction criterion, which tolerates certain pixel-level errors and evaluates the usefulness of the reconstructed image in the sense that it better serves for the downstream semantic analysis task [29]. The semantic similarity metric is successfully applied in low bitrate facial image compression [9], [46], where aggressive compression ratio is realized by filtering out the task irrelevant information. However, their supported analysis task is limited to face recognition and cannot be generalized to other applications. A task-driven semantic coding with the traditional hybrid coding framework integrates the semantic fidelity metric into the optimization process, and implements the semantic bit allocation based on reinforcement learning [28]. However, its traditional HEVC framework still process the image at pixel level and lacks the interpretation ability. A detection-driven image compression with semantically structured bit-stream is proposed in [19], where each part of the bitstream represents a specific object. Later on, the work [42] generalizes the semantically structured image coding framework to multiple intelligent tasks. However, it follows the common practice in the traditional hybrid compression framework and utilize predictive coding to reconstruct the image in pixel-level. Torfason et al. [46] verified

that semantic analysis such as classification and segmentation, can be performed on the bit-stream directly without image decoding. However, it cannot produce high quality natural image reconstruction, and the bitstream still describes the entire image without semantic structure.

Moreover, all of the above methods cannot meet the semantic and perception [4] performance requirements jointly in low bit rate. Actually, different regions of the image vary in semantic importance and should be encoded adaptively with appropriate bitrate. Reinforcement learning is a promising method to extract the task-related information and adaptively encode different regions by selecting the optimal discrete quantization coefficients. Moreover, the adversarial loss of GANs [16], [32] captures the distribution of natural images and coincide with human perception. However, it is not feasible to directly adopt the naive GANs for image semantic decoding, since the generated image tend to deviate significantly from the input image in fine-grained features [2], [3], [39].

In this paper, we propose an RL-based adaptive semantic coding (RL-ASC) approach to jointly address the semantic similarity and perceptual performance issues for ISC. The bandwidth burden could be remarkably alleviated due to the relaxing of precision requirements. Our contribution are summarized as follows:

- We propose a novel representation unit named *semantic concept* that contains the category, feature and spatial relation of each object. A convolutional semantic encoder is proposed to extract semantic concepts at the transmitter and transform them into bit streams.

- We propose a task-driven image semantic coding framework that adopts new reconstruction criterion *rate-semantic-perceptual* loss. We present an RL-based semantic bit allocation model to assign the optimal discrete quantization coefficients and achieve adaptive coding.

- We propose an attention-based generative semantic decoder at the receiver to reconstruct the image with high perceptual quality in an aggressive compression ratio. Particularly, both global and local generators are learned to capture the global and category-wise features.

The rest of this work is organized as follows. In Section II, we propose the RL-ASC method to optimize the rate-semantic-perceptual criterion. In Section III, the soft quantizer, the network architecture of attention-based generative semantic decoder as well as the training algorithm are proposed. The experiment details and the performance are shown in Section IV, and the conclusions are made in Section V.
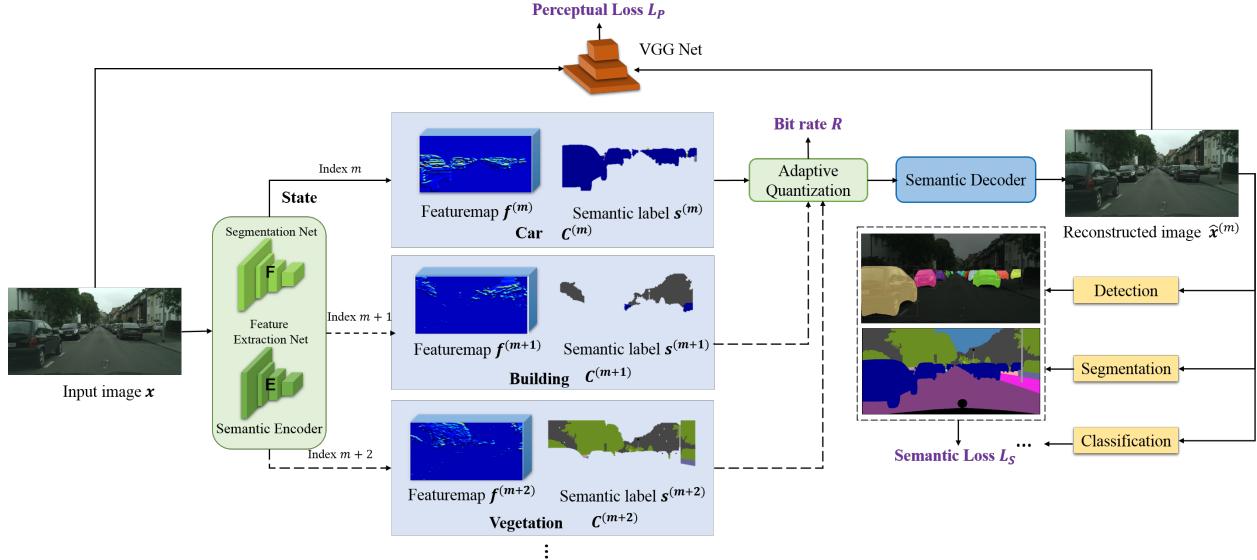
**Fig. 1:** System model of the proposed RL-based adaptive semantic coding method.

## II. THE RL-BASED ADAPTIVE SEMANTIC CODING

In this section, we propose the semantic concept as the representation unit and adopt semantic fidelity and perceptual quality as the optimization criterion. Next, we design an RL-based semantic bit allocation model to highlight and encode the task-related semantic concepts adaptively.

The system model of the proposed RL-ASC is illustrated in Fig. 1, which mainly includes three modules: the semantic encoder, the agent $\pi$ that conducts adaptive quantization, and the semantic decoder. In the inference stage, the three modules are pretrained and fixed to accommodating to a certain downstream task. The semantic encoder extracts the semantic information of each class of the given input image. Then, each class of objects is assigned with an optimal quantization level learned by the RL agent $\pi$ to assure aggressive bit rate $R$. At the receiver, the semantic decoder reconstructs all the semantic information in parallel that performs well in both semantic loss $L_s$ and perceptual loss $L_p$.

### A. Representation Unit: From Pixel to Semantic Concept

Existing image communications takes pixel as the representation unit and forces the decoded image to appear exactly the same as the transmitted one. On the contrary, the ISC interprets the underlying semantics of the image and thus is tolerant to certain pixel errors. Here we consider

the urban street scene scenario and adopt the *Cityscapes* dataset, where the quantity and classes of the semantic concepts are predefined. The dataset contains $M = 30$ number of classes grouped into the following categories: flat surfaces, humans, vehicles, constructions, objects, nature, sky, and void. If the image details cannot fit into any classes, then it belongs to the void category.

Denote the input image as $\boldsymbol{x} \in \mathbb{R}^{3 \times H \times W}$, where "3" denotes the RGB channels and $H$, $W$ are the dimensions along the height and width directions. We design the semantic encoder at the transmitter to extract semantic information for individual categories. The input image $\boldsymbol{x}$ is fed into the segmentation network $\mathcal{F}$ to obtain semantic label map: $\boldsymbol{s} = \mathcal{F}(\boldsymbol{x})$, where $\boldsymbol{s} \in \mathbb{R}^{W \times H}$ assigns a class label to each pixel in $\boldsymbol{x}$ and thus reveals the spatial arrangement of the total $M$ objects. The entry of $\boldsymbol{s}$ at pixel coordinate $(i, j)$ is an integer $\boldsymbol{s}(i, j) \in \{1, 2, ..., M\}$ that represents the class label. We denote $\boldsymbol{s}^{(m)}$ where $m = 1, 2, ..., M$ is the label ID, as the *semantic mask* for the $m$th semantic concept $\mathbf{C}^{(m)}$, the entry of which is expressed as:

$$\boldsymbol{s}^{(m)}(i, j) = \begin{cases} 1, & \text{if } \boldsymbol{s}(i, j) = m, \\ 0, & \text{otherwise.} \end{cases} \tag{1}$$

To meet the transmission bit rate requirements, the semantic information should be down-scaled and then transform into bit stream. The down-scaled semantic mask is denoted as $\boldsymbol{s}_d^{(m)} \in \mathbb{R}^{w \times h}$, where $w \leq W, h \leq H$.

Meanwhile, $\boldsymbol{x}$ is fed into the feature extraction network $\mathcal{E}$ that convolutionally processes $\boldsymbol{x}$ into a feature map of size $N \times w \times h$. Then, this feature map is projected down to $n$ channels ($n < N$) to downscale the dimension and discard perceptually redundant features, in order to meet the bit rate requirements, resulting in a feature map $\boldsymbol{f} \in \mathbb{R}^{n \times w \times h}$. Note that each convolutional layer is followed by instance normalization and LeakyReLU activation. The category-specific feature maps $\boldsymbol{f}^{(m)} \in \mathbb{R}^{n \times w \times h}$ for $\mathbf{C}^{(m)}$ is obtained by decomposing $\boldsymbol{f}$:

$$\boldsymbol{f}^{(m)} = \boldsymbol{f} \odot \boldsymbol{s}_d^{(m)}, \tag{2}$$

where $\odot$ is the element-wise multiplication.

The filtered category-specific feature maps $\boldsymbol{f}^{(m)}$ should be discriminative enough. We design a novel classification module that adopts the feature classification loss $L_C$ to assure the distinctiveness. The architecture of the classification module is as follows. Firstly, all the $\boldsymbol{f}^{(m)}$'s are fed into a max-pooling layer to yield $M$ pooled feature maps with dimension of $n \times 1 \times 1$.

Next, an fully connected (FC) layer with parameters shared across different $\boldsymbol{f}^{(m)}$'s is adopted, and the outputs are $M$ logits with dimension $M \times 1$. The subsequent softmax function returns the predicted classification probability $\hat{p}^{(m)}$ for each category-specific feature maps $\boldsymbol{f}^{(m)}$, while the one-hot ground truth label is denoted as $p^{(m)}$. The feature classification loss $L_C$ can then be written as the cross entropy loss

$$L_C = -\sum_{m=1}^{M} p^{(m)} \log \hat{p}^{(m)}. \tag{3}$$

In such a way, we imitate the image understanding process of humans and divide $\boldsymbol{x}$ into $M$ individual semantic concepts that each can be denoted as a set $\mathbf{C}^{(m)} = \{\boldsymbol{f}^{(m)} \in \mathbb{R}^{n \times w \times h}, \boldsymbol{s}_d^{(m)} \in \mathbb{R}^{w \times h}\}$. Therefore, $\mathbf{C}^{(m)}$ can be deemed as novel representation unit that is far more efficient than pixel unit.

## B. Reconstruction Metric: From Pixel Loss to Semantic-Perceptual Loss

Conventional image communications system measures the pixel loss in terms of PSNR and the coding process is optimized by the rate-distortion theory. However, the main purpose of ISC is to transmit the underlying meaning of the image other than pixel copy. The ISC should extract, encode, and reconstruct each $\mathbf{C}^{(m)}$ with low bitrate in the sense that critical concepts are recovered with high semantic fidelity and perceptual quality. We propose a novel rate-semantic-perceptual criterion to optimize the semantic coding process.

It is already known that the intelligent analysis tasks such as object detection, semantic segmentation, pose estimation, and action recognition, etc., are able to extract semantic information accurately, and the obtained semantic information could also be used in the image coding process. Denote the network of the intelligent task as $\mathcal{H}$, and then the prediction result of the input image $\boldsymbol{x}$ and reconstructed image $\hat{\boldsymbol{x}}$ can be written as $\mathcal{H}(\boldsymbol{x})$ and $\mathcal{H}(\hat{\boldsymbol{x}})$. We propose a novel semantic loss $L_S$ metric defined as the degradation of the precision performance of the intelligent task $\mathcal{H}$ on the reconstructed image $\hat{\boldsymbol{x}}$ compared to the original image $\boldsymbol{x}$. The $L_S$ can even be generalized to arbitrary user-defined semantic tasks, which can be written as:

$$L_S = f_{degrade}(\mathcal{H}(\boldsymbol{x}), \mathcal{H}(\hat{\boldsymbol{x}})), \tag{4}$$

where $f_{degrade}$ is the accuracy degradation function of the predicted result given the ground truth. In such a way, we successfully convert abstract semantic information to a measurable form.

The PSPNet [53] can be leveraged to implement the semantic segmentation task $\mathcal{H}$. while the bounding set of the predicted object of class $m$ on $\hat{\boldsymbol{x}}$ is $\mathcal{H}(\hat{\boldsymbol{x}}) = \boldsymbol{B}_d^m$, while the ground truth bounding set on $\boldsymbol{x}$ is denoted as $\mathcal{H}(\boldsymbol{x}) = \boldsymbol{B}_g^m$. The performance of semantic segmentation can be evaluated by IoU that measures the consistency of $\boldsymbol{B}_g^m$ and $\boldsymbol{B}_d^m$:

$$IoU(\mathcal{H}(\boldsymbol{x}), \mathcal{H}(\hat{\boldsymbol{x}})) = IoU(\boldsymbol{B}_g^m, \boldsymbol{B}_d^m) = overlap(\boldsymbol{B}_g^m, \boldsymbol{B}_d^m)/union(\boldsymbol{B}_g^m, \boldsymbol{B}_d^m), \tag{5}$$

where IoU is a value between 0 and 1. Mean IoU (mIoU) is the average segmentation precision over the total $M$ classes. The semantic loss is calculated by mIoU and represents the extent of the non-interpretability of the reconstructed image $\hat{\boldsymbol{x}}$:

$$L_S = 1 - \frac{1}{M} \sum_{m=1}^{M} IoU(\boldsymbol{B}_g^m, \boldsymbol{B}_d^m). \tag{6}$$

The ideal $\hat{\boldsymbol{x}}$ should preserve certain semantic information required by PSPNet, which results in $L_s = 0$.

The Mask R-CNN [18] can be utilized to implement the object detection task $\mathcal{H}$. We adopt IoU loss [52] as the semantic loss:

$$L_S = -\ln[\frac{1}{M} \sum_{m=1}^{M} IoU(\boldsymbol{B}_g^m, \boldsymbol{B}_d^m)]. \tag{7}$$

For image classification task, the semantic loss $L_S$ is defined as the typical cross-entropy function. Let $\mathcal{H}$ be the pretrained classification network such as VGG and MobileNet [38]. Let $\mathcal{H}(\boldsymbol{x}) = p_1, p_2, ..., p_M$ be the one-hot vector of the ground truth label. Denote $y$ as the ground truth label ID, and therefore we can obtain

$$p_i = \begin{cases} 1, & if(i = y) \\ 0, & if(i \neq y) \end{cases} \tag{8}$$

The predicted output is denoted as $\hat{\boldsymbol{x}} = \hat{p}_1, \hat{p}_2, ..., \hat{p}_M$ corresponding to $M$ classes. The semantic loss is computed as

$$L_S = -\frac{1}{M} \sum_{i=1}^{M} (p_i \log(\hat{p}_i)), \tag{9}$$

For person re-identification, the performance can be evaluated by the cross-entropy loss (9) utilizing the label smoothing technique. The true probability distribution is rewritten as

$$p_i = \begin{cases} 1 - \epsilon, & if(i = y) \\ \epsilon/(K-1), & if(i \neq y) \end{cases} \tag{10}$$

where $\epsilon$ is a small hyper parameter.

Besides, the visual performance such as the naturalness and clarity of $\hat{\boldsymbol{x}}$ is another critical issue for the ISC. The perceptual loss projects the images $\boldsymbol{x}$ and $\hat{\boldsymbol{x}}$ to a high dimensional feature space with a pretrained model and minimize the distance in the high-level space to measure the similarity of the two images. The perceptual loss is validated to predict human scores for degradation properly and could be expressed as:

$$L_P = \|\phi(\boldsymbol{x}) - \phi(\hat{\boldsymbol{x}})\|_2^2, \tag{11}$$

where $\phi(\cdot)$ is the pretrained VGG network that maps an input image to a high dimension feature space.

We then formulate the rate-semantic-perceptual criterion $L$ as the triple trade-off between semantic fidelity $L_S$, perceptual quality $L_P$, and bit rates $R$ of the semantic concepts:

$$L = \lambda R + L_S + \eta L_P \tag{12}$$

where $R$ is measured by the length of transmitted bit stream, and $\lambda, \eta$ are the weighting parameters that balance the three terms. The different points of the rate-semantic-perceptual curve can be obtained by changing the value of $\lambda$.

## C. The RL-based Semantic Bit Allocation Model

Different semantic concepts matter differently in downstream analysis tasks, and only the task-related semantic information needs to be transmitted with high precision. The critical semantic concepts should be focused on and then encoded precisely without semantic loss, while the precision requirements for other objects can be relaxed. Taking the street scene dataset Cityscapes as an example, the object detection task mainly concentrates on the vehicle, pedestrian, and traffic sign, etc., while the background such as sky and vegetation does not affect the core meaning of the image. Besides, a few bits can already well represent sky and vegetation, since their texture is simple and regular.

In order to locate the critical semantic concepts and assign appropriate precision requirement, we propose an RL-based semantic bit allocation model. After all the semantic concepts are obtained by the semantic encoder, the coding process is modeled as a Markov Decision Process
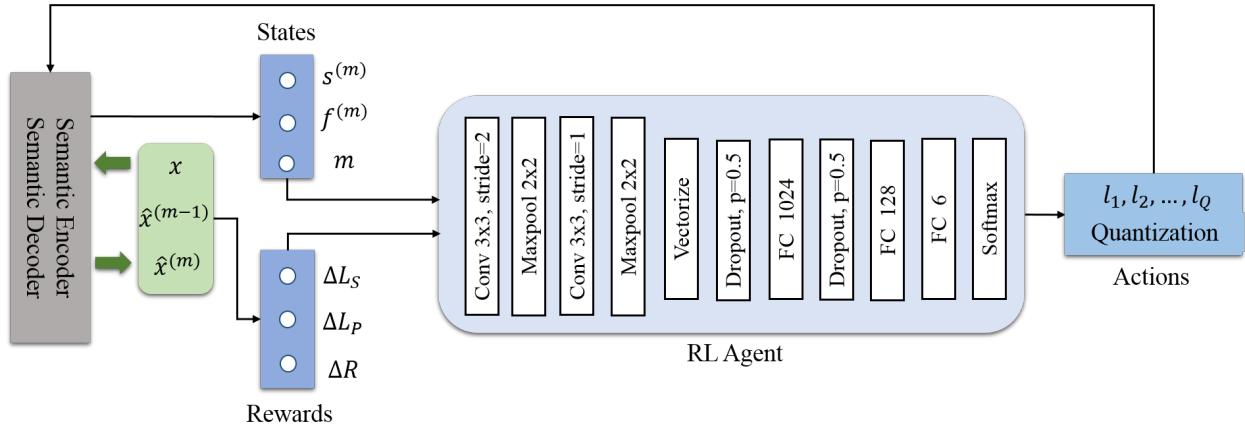
**Fig. 2:** The architecture of the RL-based Semantic Bit Allocation Model.

(MDP), where each semantic concept is encoded and decoded sequentially with appropriate bit rate in the label index order from $m = 1$ to $m = M$. The RL agent adaptively assigns the quantization level for each semantic concept by balancing the bit cost and reconstruction performance. Accordingly, we can integrated the semantic fidelity and perceptual quality metrics into the semantic coding optimization. Therefore, the more informative and visually salient objects are represented with higher bitrates while the irrelevant regions are represented with lower bitrates. When the training of the RL agent is completed, the whole process of the quantization decision is off-policy, which can be processed in parallel with the encoder.

The RL framework typically consists of a state space $\mathbb{S}$ that includes all the possible information that the agent can observe from the environment. The state in the $m$th step is $state^{(m)} = \{\boldsymbol{f}^{(m)}, \boldsymbol{s}^{(m)}, m\}$, and we denote $\hat{\boldsymbol{x}}^{(m-1)}$ as the reconstructed image from the former step. The action space $\mathbb{A}$ is defined as a set of quantization levels $action^{(m)} \in \{l_1, l_2, ..., l_Q\}$ assigned for the semantic concepts, where $Q$ is the total number of levels and a higher quantization level leads to a higher bit rate and finer details. We define the behavior of the agent as a policy $\pi : \mathbb{S} \times \mathbb{A} \rightarrow [0, 1]$ that maps states to a distribution of actions, denoted as $\pi(action^{(m)}|state^{(m)})$. Specifically, $\pi$ can be implemented by a neural network parameterized by $\theta$. The intermediate reward $r^{(m+1)}$ is received after the agent $\pi$ taking action $action^{(m)}$ at sate $state^{(m)}$ by evaluating the performance of the current reconstructed image $\hat{\boldsymbol{x}}^{(m)}$. The transition function $trans(state^{(m)}, action^{(m)})$ pre-

dicts the next state given the current state and the action. In this work, the evolution from state $state^{(m)}$ to $state^{(m+1)}$ is deterministic, and thus $p(state^{(m+1)}) = trans(state^{(m)}, action^{(m)}) = 1$.

As shown in Fig. 2, the detailed procedure of the agent $\pi$ is as follows. The current state $state^{(m)}$ is fed into the agent to learn the policy $\pi$ that reveals the semantic importance and accordingly produces an action $action^{(m)}$ that adaptively quantizes $\boldsymbol{C}^{(m)}$. Next, the semantic decoder takes the received bitstream and the selected quantization level as input to produce the current reconstructed image $\hat{\boldsymbol{x}}^{(m)}$. The immediate reward $r^{(m+1)}$ is defined as the quality increase from $\hat{\boldsymbol{x}}^{(m)}$ to $\hat{\boldsymbol{x}}^{(m+1)}$ in terms of the increase of semantic similarity, perceptual quality, as well as the decrease in bitrate $R$. Note that compared to $\hat{\boldsymbol{x}}^{(m-1)}$, the current $\hat{\boldsymbol{x}}^{(m)}$ is updated in the region of the particular semantic concept $\boldsymbol{C}^{(m)}$, while other region maintains unchanged. The bitrate $R$ of $\hat{\boldsymbol{x}}^{(m)}$ is denoted as $\psi(\hat{\boldsymbol{x}}^{(m)})$ and computed as the bitrate summation of the so far processed category-specific features $\boldsymbol{f}^{(m)}$'s as well as the segmentation map $\boldsymbol{s}$. The detailed bitrate calculation method is expressed in Section III. A.

We next design the architecture of RL agent $\pi$. Firstly, the state $state^{(m)}$ is fed into two consecutive convolutional layers, each of which is followed by the ReLU activation and max-pooling operation. Next, the output of the convolutional layer is flattened into a vector and then fed into three FC layers to yield an $M \times 1$ dimensional logits. Note that the dropout layer and ReLU activation are adopted to avoid over-fitting. Thirdly, the subsequent softmax function returns the predicted probability of each quantization level, and the action $action^{(m)}$ is sampled from the predicted probability distribution.

At the initialization step, we set the coarsest quantization level $l_1$ for the whole image to obtain an initial reconstruction $\hat{\boldsymbol{x}}^{(0)}$ with the lowest bit rate $\psi(\hat{\boldsymbol{x}}^{(0)})$. Then, both of the two images $\boldsymbol{x}$ and $\hat{\boldsymbol{x}}^{(0)}$ are fed into the downstream analysis network $\mathcal{H}$ (such as PSPNet) and the VGG network to measure the semantic loss and perceptual loss. At the next step, the state $state^{(1)} = \{\boldsymbol{f}^{(1)}, \boldsymbol{s}^{(1)}, m = 1\}$ is take into consideration, and the action $action^{(1)}$ is produced by the agent to select a proper quantization level. The current reconstructed image $\hat{\boldsymbol{x}}^{(1)}$ is obtained by the semantic decoder and sequentially the reconstruction performance is evaluated. At step $m$, the rate-semantic-perceptual loss can be written as

$$L^{(m)} = \lambda[\psi(\hat{\boldsymbol{x}}^{(m)}) - \psi(\hat{\boldsymbol{x}}^{(0)})] + f_{degrade}(\mathcal{H}(\boldsymbol{x}), \mathcal{H}(\hat{\boldsymbol{x}}^{(m)})) + \eta\|\phi(\boldsymbol{x}) - \phi(\hat{\boldsymbol{x}}^{(m)})\|_2. \quad (13)$$

At step $m+1$, the $(m+1)$th semantic concept $\mathbf{C}^{(m+1)}$ on the current reconstruction $\hat{\boldsymbol{x}}^{(m+1)}$ is updated and we will accordingly obtain a renewed rate-semantic-perceptual loss $L^{(m+1)}$. The intermediate reward could be written as

$$r^{(m+1)} = L^{(m)} - L^{(m+1)} = \lambda \triangle R + \triangle L_S + \eta \triangle L_P, \tag{14}$$

where $\triangle L_S, \triangle L_P$ and $\triangle R$ represent the difference in the semantic loss, the perceptual loss, and the bit rate before and after coding.

The agent should be driven to maximize the cumulated reward in the whole episode, which is the optimal goal of the MDP. Denote $\gamma$ as the discounted factor, the discounted cumulated reward for a full episode can be written as

$$G = \sum_{m=1}^{M} \gamma^m r^{(m+1)}. \tag{15}$$

The objective function for such a learning process can be formulated as maximizing the expected $G$ for all trajectories. Denote $T$ as a complete trajectory $(state^{(1)}, action^{(1)}, r^{(2)}, state^{(2)}, action^{(2)}, r^{(3)}, ..., state^{(M)}, action^{(M)}, r^{(M+1)})$, and $T_\pi$ as the trajectory distribution. Denote the total number of sampled trajectories as $N$. Then the objective function can be written as

$$J(\pi) = \mathbb{E}_{T_\pi} G = \sum_N T_\pi G. \tag{16}$$

The optimal policy $\pi$ can be obtained by performing gradient ascend method on the sampled trajectories. The derivative of $J_\pi$ can be calculated as

$$\nabla_\theta J_\pi = \sum_N G \nabla_\theta T_\pi = \sum_N G T_\pi \frac{\nabla_\theta T_\pi}{T_\pi} = \sum_N T_\pi G \nabla_\theta \log T_\pi. \tag{17}$$

Since $trans(state^{(m)}, action^{(m)})$ is deterministic, we may further expand $T_\pi$ as

$$
\begin{aligned}
T_\pi &= p(state^{(1)}) \prod_{m=1}^{M} \pi(action^{(m)}|state^{(m)}) trans(state^{(m)}, action^{(m)}) \\
&= p(state^{(1)}) \prod_{m=1}^{M} \pi(action^{(m)}|state^{(m)}),
\end{aligned}
\tag{18}
$$

where $p(state^{(1)})$ is the probability of observing $state^{(1)}$. For computational efficiency, we calculate the derivative of the logarithm form of formula (18) as

$$
\begin{aligned}
\nabla_\theta \log T_\pi &= \nabla_\theta \log p(state^{(1)}) + \sum_{m=1}^{M} \nabla_\theta \log \pi(action^{(m)}|state^{(m)}) \\
&= \sum_{m=1}^{M} \nabla_\theta \log \pi(action^{(m)}|state^{(m)}),
\end{aligned}
\tag{19}
$$

where the first term $\log p(state^{(1)})$ is irrelevant to $\theta$ and thus the derivative is zero. According to (19), equation (17) can be rewritten as

$$\nabla_\theta J_\pi = \mathbb{E}_{T_\pi}[G \sum_{m=1}^{M} \nabla_\theta \log \pi(action^{(m)}|state^{(m)})]. \tag{20}$$

In practice, we can use a one-time Monte-Carlo rollout to sample a trajectory. Therefore, equation (20) can be written as:

$$\nabla_\theta J_\pi = G \sum_{m=1}^{M} \nabla_\theta \log \pi(action^{(m)}|state^{(m)}). \tag{21}$$

Then, the parameters of the RL agent $\pi$ is updated as

$$\theta \leftarrow \theta + \alpha \nabla_\theta J_\pi, \tag{22}$$

where $\alpha$ is the learning rate.

## III. SEMANTIC DECODER AND TRAINING DETAILS

As shown in Fig. 3, we design the soft quantization, entropy coding, semantic decoder and the training details in this section.

### A. Soft Quantization and Entropy Coding

The RL agent $\pi$ selects the action $a^{(m)}$ for each $\boldsymbol{f}^{(m)}$ to assign bits adaptively and concentrate on semantic important regions. We adopt the scalar variant for the quantization approach and obtain the discrete $\hat{\boldsymbol{f}}^{(m)}$ as

$$\hat{\boldsymbol{f}}^{(m)} := Quantize(\boldsymbol{f}^{(m)}, a^{(m)}), \tag{23}$$

where each entry of $\hat{\boldsymbol{f}}_m$ at coordinate $(k, i, j)$ can be computed by the nearest neighbor assignment method. The process can be written as

$$\hat{\boldsymbol{f}}_m(k, i, j) = \arg\min_t \|\boldsymbol{f}^{(m)}(k, i, j) - l_m\|, \tag{24}$$

where $l_m \in \{l_1, l_2, ..., l_Q\}$ is the quantization center.

A major problem in quantization is that the gradients of (24) are zeros almost everywhere, which makes gradient descent-based optimization ineffective in the end-to-end communication
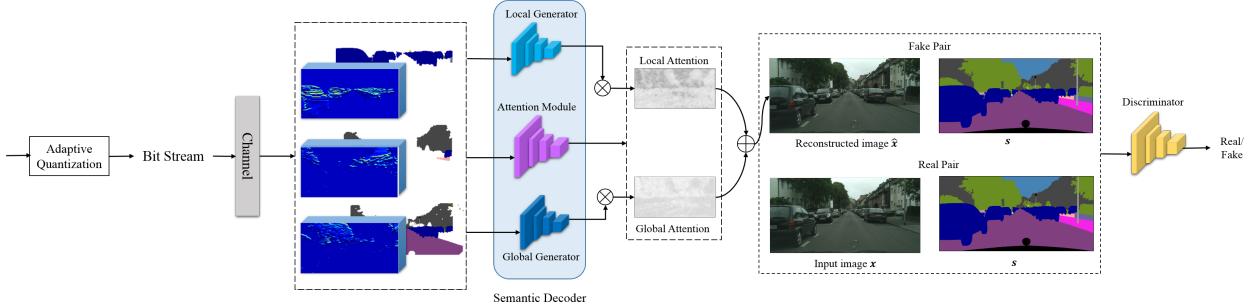
**Fig. 3:** The architecture of the semantic decoder.

system. To address the zero gradient problem and approximate formulation (24), we introduce a differentiable soft quantization with continuous functions [15]

$$\tilde{\boldsymbol{f}}^{(m)}(i,j,k) = \sum_{m=1}^{Q} \frac{exp(-\sigma\|\boldsymbol{f}^{(m)}(k,i,j) - l_m\|)}{\sum_{t=1}^{Q} exp(-\sigma\|\boldsymbol{f}^{(m)}(k,i,j) - l_t\|)} l_m,$$ (25)

where $\sigma$ is a hyperparameter relating to the softness of the quantization. Note that the quantization formulation (24) is adopted in the forward pass, while the differentiable soft quantization formulation (25) is adopted in the backward pass.

It is already known that the entropy coding can reduce the expected bit length by assigning a short code for frequently occurred codeword while assigning a long code for others. We adopt Huffman coding as the entropy coding, where the codeword probability is predicted by building and maintaining a frequency table. In this case, $\hat{\boldsymbol{f}}^{(m)}$ is transformed into the variable-length binary code $\boldsymbol{r}^{(m)}$ of length $\ell(\boldsymbol{r}^{(m)})$ as

$$\boldsymbol{r}^{(m)} := Huffman(\hat{\boldsymbol{f}}^{(m)}) \in \{0,1\}^{\ell(\boldsymbol{r}^{(m)})}.$$ (26)

Besides, the semantic label map $\boldsymbol{s}$ is losslessly coded in vector graphic format, and the binary code is denoted as $\boldsymbol{r}_s$. The coded bitstream of the proposed RL-ASC includes two parts: (i) the binary codes $\boldsymbol{r}_s$ with fixed length $\ell(\boldsymbol{r_s})$; (ii) the binary codes $\boldsymbol{r}^{(m)}$ with variable length $\ell(\boldsymbol{r}_m)$ that is adjusted by $a^{(m)}$. We thereby define the rate of the reconstructed image $\hat{\boldsymbol{x}}^{(m)}$ at step $m$ as the total bits of all the binary codes $\boldsymbol{r}^{(m)}$ and $\boldsymbol{r}_s$, which can be expressed as

$$\psi(\hat{\boldsymbol{x}}^{(m)}) = \ell(\boldsymbol{r_s}) + \sum_{m=1}^{M} (\ell(\boldsymbol{r}^{(m)})).$$ (27)

Then, the binary code $\boldsymbol{r}^{(m)}$ and $\boldsymbol{r}_s$ are transmitted through the wireless channel and arrive at the receiver. Unless notified, we assume the wireless channel is lossless since the paper
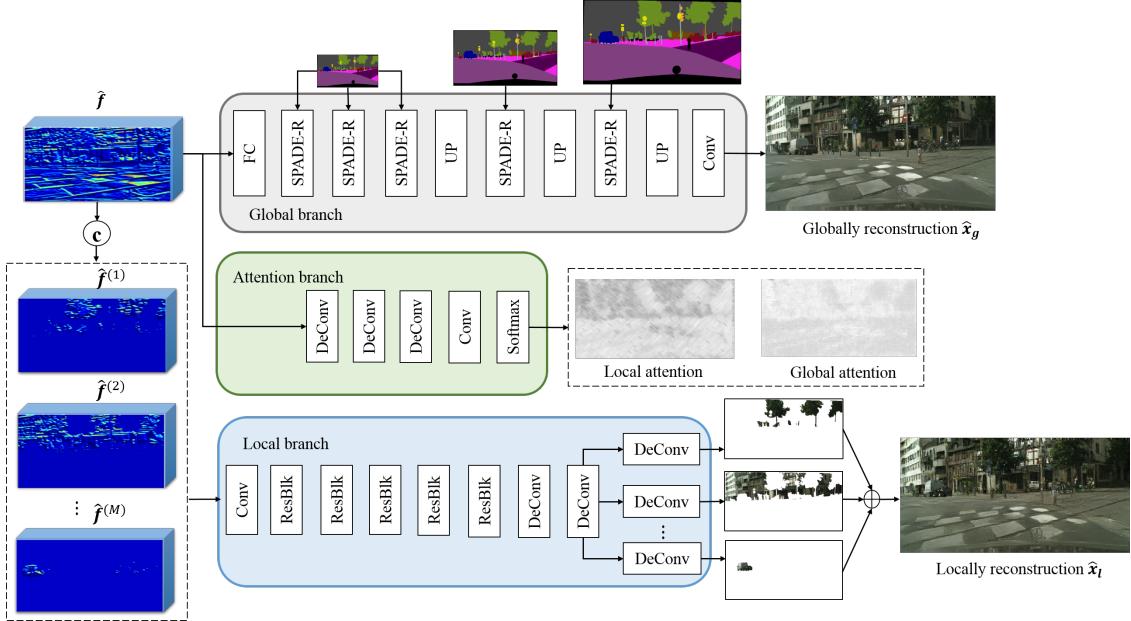
**Fig. 4:** The architecture of the generative semantic decoder.

mainly focuses on image semantic encoding and semantic decoding. Therefore, the receiver could achieve perfect $\hat{\boldsymbol{f}}^{(m)}$ and $\boldsymbol{s}$. Benefiting from the incredible semantic extraction and bit allocation performance, the essential semantics are appropriately preserved and transmitted.

### B. Generative Semantic Decoder

In this section, the received bit streams are decoded into semantic concepts and then combined as the reconstructed image. Some learning-based codecs are based on convolutional neural networks (CNNs) [4], [21], [26], [33], while others on recurrent neural networks (RNNs) [23], [45]. Recent learning-based codecs [10], [41] are competitive with or even superior to the classical codecs. The work [10] proposed to use discretized Gaussian Mixture Likelihoods to parameterize the distributions of latent codes, which can achieve a more accurate and flexible entropy model. The work [41] propose a versatile deep image compression network based on Spatial Feature Transform, which enables task-aware image compression for various tasks, e.g., classification, with variable rates. However, all of the aforementioned works optimizes the rate-distortion (RD) performance and cannot meet the requirements of ISC.

Different from conventional methods that reduce pixel errors, the proposed semantic decoder should be tolerant of pixel errors to some extent. GANs [16] has shown impressive ability in

synthesizing high-resolution images. Instead of forcing the generated images to be exactly the same as the ground truth image, GANs attempt to approximate the intractable distributions of the real image dataset. Therefore, the generated image maintains the underlying meaning of the ground truth image while the two images may differ in pixel values. Therefore, GANs fits the spirit of the ISC and are adopted as the semantic decoder in the proposed method.

Adversarial training of the GANs has shown advantage in reducing compression artifacts in deep compression systems [3], [4], [37], [43]. Inspired by existing methods, we adopt conditional GANs in the semantic decoder that translates the semantic concepts into the reconstructed image. The reconstructed image should be natural and similar to the input image. The naturalness can be measured by the adversarial loss that indicates the probability distributions divergence of real and reconstructed images, while the visual similarity could be measured by the aforementioned perceptual loss that indicates the distance in feature space [7].

The detailed architecture of the generative semantic decoder is illustrated in Fig. 4, which consists of a local generator $\mathcal{G}_l$, a global generator $\mathcal{G}_g$, and an attention module.

**Local Generator** $\mathcal{G}_l$**:** The quantity and spatial occupation of different categories are imbalanced in the training dataset. It is extremely difficult to generate small object classes and texture details since the dataset are dominated by frequently occurred or large object classes. To avoid the interference from other classes, we propose a novel local image generator $\mathcal{G}^l$ to separately reconstruct the class-specific objects $\boldsymbol{y}^{(m)} \in \mathbb{R}^{H \times W}$ from the quantized class-specific feature maps $\hat{\boldsymbol{f}}^{(m)}$ as

$$\boldsymbol{y}^{(m)} = \mathcal{G}^l(\hat{\boldsymbol{f}}^{(m)}). \tag{28}$$

The locally reconstructed image $\hat{\boldsymbol{x}}_l \in \mathbb{R}^{H \times W}$ is obtained by element-wise addition of all the individual $\boldsymbol{y}^{(m)}$'s as

$$\hat{\boldsymbol{x}}_l = \boldsymbol{y}^{(1)} \oplus \boldsymbol{y}^{(2)} \oplus \cdots \oplus \boldsymbol{y}^{(M)}. \tag{29}$$

As shown in Fig. 4, we feed $M$ feature maps $\hat{\boldsymbol{f}}^{(m)}$ into a convolution layer and five consecutive residual blocks, and the outputs are then up-scaled by two consecutive deconvolutional layers. Then, the up-scaled feature maps of $M$ semantic concepts are separately fed into $M$ corresponding deconvolutional layers to produce $\boldsymbol{y}^{(m)}$. Each deconvolutional layer has independent network parameters and is able to effectively preserve the class-specific features with rich details.

**Global Generator** $\mathcal{G}_g$**:** We attempt to capture the global structure information as well as the spatial layout by global generation. Inspired from GauGAN [34], the well designed $\mathcal{G}_g$ adopts spatially adaptive normalization (SPADE) to fuse the spatial layout information from $s$ into the feature maps $\hat{f}$. The SPADE modulates the layer activation in accordance with $s$ to guide the reconstruction of different categories. The globally reconstructed image $\hat{x}_g$ can be obtained as

$$\hat{x}_g = \mathcal{G}_g(\hat{f}, s). \tag{30}$$

Particularly, we feed $\hat{f}$ into three residual blocks with SPADE (SPADE-R), the output of which is upsampled three times to yield the global reconstruction $\hat{x}_g$.

**Attention Module:** In order to better combine the outputs of $\mathcal{G}_l$ and $\mathcal{G}_g$, we further propose an attention module to learn the local weight matrix $W_l \in \mathbb{R}^{H \times W}$ and global weight matrix $W_g \in \mathbb{R}^{H \times W}$. The feature map $\hat{f}$ is fed into the attention module that includes three consecutive deconvolutional layers and a convolutional layer. The output of the convolutional layer has two channels that correspond to the two weight matrices and is further normalized by a softmax layer in channel-wise. Therefore, we obtain $W_g$ and $W_l$ whose values are in the range of (0,1) and meet the condition $W_l(i,j) + W_g(i,j) = 1$ in each location $(i,j)$. The final output image $\hat{x}$ can be written as element-wise summation of weighted local reconstruction $\hat{x}_l$ and global reconstruction $\hat{x}_g$:

$$\hat{x} = (W_l \odot \hat{x}_l) \oplus (W_g \odot \hat{x}_g). \tag{31}$$

We adopt a multi-scale patch-discriminator $\mathcal{D}$ to identify the real image $x$ and the final output image $\hat{x}$. The generator $\mathcal{G}_l, \mathcal{G}_g$ are trained alternatively with the discriminator $\mathcal{D}$. Denote $\mathcal{D}(x, s)$ or $\mathcal{D}(\hat{x}, s)$ as the predicted probability that the image pair $(x, s)$ or $(\hat{x}, s)$ are from real samples. The discriminator $\mathcal{D}$ learns to distinguish real image pairs from fake ones by maximizing $\mathcal{D}(x, s)$ and minimizing $\mathcal{D}(\hat{x}, s)$. The generators $\mathcal{G}_l$ and $\mathcal{G}_g$ learn to fool $\mathcal{D}$ by minimizing $(-\mathcal{D}(\hat{x}, s))$. We adopt the hinge loss for GANs, whose objective function is

$$\min_{\mathcal{D}} L_{\mathcal{D}} = \mathbb{E}[\max(0, 1 - \mathcal{D}(x, s))] + \mathbb{E}[\max(0, 1 + \mathcal{D}(\hat{x}, s))],$$

$$\min_{\mathcal{G}_l, \mathcal{G}_g} L_{\mathcal{G}} = -\mathbb{E}[\mathcal{D}(\hat{x}, s)]. \tag{32}$$

Note that the hinge loss $L_{\mathcal{D}}$ is used for "maximum-margin" classification, which penalizes the positive samples of $\mathcal{D}(x, s) < 1$ and the negative samples of $\mathcal{D}(\hat{x}, s) > -1$. Only the samples

that are not classified properly will impact the gradient update of $\mathcal{G}_l$, $\mathcal{G}_g$ and $\mathcal{D}$. The formulation (32) forces the probability distribution of $\hat{\boldsymbol{x}}$ to approximate that of the training image set. If the feature extraction network $\mathcal{E}$ cannot afford to store the exact detail in $\boldsymbol{x}$, then $\mathcal{G}^l$ and $\mathcal{G}^g$ are able to synthesize the detail to satisfy natural image distribution instead of showing blocky and blurry effects. Equilibrium will be achieved when $\mathcal{D}$ classifies the reconstructed image $\hat{\boldsymbol{x}}$ as real.

The final reconstruction $\hat{\boldsymbol{x}}$ should not only be consistent with the natural image distribution, which coincides with the target of traditional GANs, but also meet the requirements of semantic communications and recover the semantic concepts of a specific input image $\boldsymbol{x}$. We adopt the perceptual loss (11) to guide the semantic reconstruction process and ensure $\hat{\boldsymbol{x}}$ approximates the feature space of $\boldsymbol{x}$.

The proposed semantic encoder and semantic decoder are trained jointly, and the objective function can be formulated as the weighted combination of adversarial loss $L_{\mathcal{D}}$, perceptual loss $L_P$, and feature classification loss $L_C$ as

$$\min_{\mathcal{E}, \mathcal{G}, \mathcal{D}} L_{\mathcal{D}} + \lambda_1 L_P + \lambda_2 L_C, \tag{33}$$

where $\lambda_1$ and $\lambda_2$ are the weighting parameters. Due to the adequately preserved semantic information, the proposed generative semantic decoder produces more incredible reconstructions compared with the leading image generation models [2], [3], [34]. Note that the work [2] performs poor in preserving appearance features, and the reconstructed image deviates apparently from the source image. The work [34] can merely generate a general image from the distribution of the training dataset without the ability to control the appearance of certain objects.

## C. Training Algorithm

We next propose a three-stage training algorithm to ensure the semantic encoder, the semantic decoder, and the RL-based semantic bit allocation model function properly. Benefit from the differentiable soft quantization, the whole system can be trained with an ADAM solver [24], where the momentum term of ADAM is set to be $\beta_1 = 0.5$ and $\beta_2 = 0.999$. The whole process is to first conduct stage I and train the semantic encoder and decoder in an end-to-end manner, while taking equal quantization for every semantic concepts with the highest quantization level. Then, with the pretrained semantic encoder and decoder being fixed, we conduct stage II to

---

**Algorithm 1:** Training the semantic encoder and semantic decoder

---

**Input:** Image dataset $\boldsymbol{X}$, pretrained $\mathcal{F}$, training epoch $E_p = 70$, learning rate $\alpha = 2e^{-4}$, parameters $\lambda_1$ and $\lambda_2$, batch size $b$, total categories of semantic concepts $M$, and fixed quantization level $a^{(M)}$.

**Output:** Neural network parameter for $\mathcal{G}_l$, $\mathcal{G}_g$, $\mathcal{D}$, and $\mathcal{E}$.

**1** **for** *epoch=1:$E_p$* **do**

**2**     Sample a batch of image $\{\boldsymbol{x}_1, ..., \boldsymbol{x}_b\} \in \boldsymbol{X}$;

**3**     Obtain $M$ feature maps $\hat{\boldsymbol{f}}_1^{(m)}, ..., \hat{\boldsymbol{f}}_b^{(m)}$ and $\{\boldsymbol{s}_1, ..., \boldsymbol{s}_b\}$ for the image batch by $\mathcal{E}$ and $\mathcal{F}$;

**4**     Quantize $M$ feature maps $\hat{\boldsymbol{f}}_1^{(m)}, ..., \hat{\boldsymbol{f}}_b^{(m)}$ with quantization level $a^{(M)}$ by formulation (25);

**5**     The bit streams are transmitted to the receiver through channel;

**6**     Obtain the final output image $\{\hat{\boldsymbol{x}}_1, ..., \hat{\boldsymbol{x}}_b\}$ by $\mathcal{G}_l$ and $\mathcal{G}_g$ by formulation (31);

**7**     Update the discriminator $\mathcal{D}$ by descending its stochastic gradient:

$\nabla_{\theta_d} \frac{1}{b} \sum_{j=1}^{b} \{[\max(0, 1 - \mathcal{D}(\boldsymbol{x}_j, \boldsymbol{s}_j))] + [\max(0, 1 + \mathcal{D}(\hat{\boldsymbol{x}}_j, \boldsymbol{s}_j))]\}$.

**8**     Update the generators $\mathcal{G}_l$, $\mathcal{G}_g$ and $\mathcal{E}$ by descending their stochastic gradient:

$\nabla_{\theta_{g,e}} \frac{1}{b} \sum_{j=1}^{b} [-\mathcal{D}(\hat{\boldsymbol{x}}_j, \boldsymbol{s}_j) + \lambda_1 L_P + \lambda_2 L_C]$.

**9** **end**

**10** Return the network $\mathcal{G}_l, \mathcal{G}_g, \mathcal{D}$, and $\mathcal{E}$.

---

train the semantic bit allocation model, a.k.a, the RL agent $\pi$. In stage III, the semantic encoder, semantic decoder, and RL agent $\pi$ are finetuned together.

**Stage I Training the semantic encoder and semantic decoder**: In this stage, the proposed RL-ASC learns to extract semantic concepts and then semantically reconstruct the image, laying the basis for training the RL agent $\pi$. We first initialize $\mathcal{E}$, $\mathcal{G}_l$, $\mathcal{G}_g$, and $\mathcal{D}$ but leave the RL agent $\pi$ out of this stage. We assume all the semantic concepts $\mathbf{C}^{(m)}$'s in the input image $\boldsymbol{x}$ are of the same importance, and the quantization precision is set to be the highest level for each concept. Hence, the lowest information loss is obtained after quantization and the semantic decoder can produce the reconstruction $\hat{\boldsymbol{x}}$ as fine as possible. The training algorithm is summarized in Algorithm 1.

Alternate training is applied between one gradient descent step on $\mathcal{D}$, then one step on $\mathcal{G}_l$, $\mathcal{G}_g$, and $\mathcal{E}$. The training takes place in a loop, where the adversarial loss $L_{\mathcal{D}}$, the perceptual loss $L_P$, and the feature classification loss $L_C$ are minimized as in (33). Firstly, we update $\mathcal{D}$ by real training images and the reconstructed images. Next, with $\mathcal{D}$ being set as non-trainable, we feed $\hat{\boldsymbol{f}}^{(m)}$ and $\boldsymbol{s}$ to both $\mathcal{G}_l$ and $\mathcal{G}_g$ to produce final output image $\hat{\boldsymbol{x}}$, and then we adopt $\mathcal{D}$

---

**Algorithm 2:** Training the RL-based semantic bit allocation model

---

**Input:** Image dataset $\boldsymbol{X}$, training epoch $E_p = 5$, learning rate $\alpha = 1e^{-5}$, parameters $\eta = 10$ and $\lambda = 1$,

batch size $b = 1$, total categories of semantic concepts $M$, quantization levels $\mathbb{A} = \{l_1, l_2, ..., l_Q\}$,

discounted factor $\gamma = 0.99$, the pre-trained semantic encoder and semantic decoder.

**Output:** Neural network parameter $\theta$.

**1 for** *epoch=1:$E_p$* **do**

**2**     Sample a batch of image data $\boldsymbol{x} \in \boldsymbol{X}$.

**3**     Initialization: Set all the $M$ semantic concepts $\mathbf{C}^{(m)}$'s with the coarsest quantization level $a^{(0)} = l_1$

      and obtain the bitrate $\psi(\boldsymbol{x}^{(0)})$ by (27).

**4**     **for** *m=1:M* **do**

**5**        **Encoding**: The semantic encoder extracts $\mathbf{C}^{(m)}$ from $\boldsymbol{x}$;

**6**        **Quantization**: The policy $\pi(a^{(m)}|s^{(m)})$ selects an action $a^{(m)} \in \mathbb{A}$ for $\mathbf{C}^{(m)}$;

**7**        **Entropy Coding**: Conduct Huffman coding by formulation (26) and obtain bitrate $\psi(\boldsymbol{x}^{(m)})$;

**8**        **Transmitting**: The bit streams are transmitted to the receiver through channel;

**9**        **Decoding**: The semantic decoder reconstructs $\hat{\boldsymbol{x}}^{(m)}$ by formulation (31);

**10**        **Reward**: Calculate intermediate reward $r^{(m+1)}$ by formulation (13).

**11**     **end**

**12**     Calculate discounted cumulated reward $G$ by formulation (15).

**13**     Calculate the gradient of $J_\theta$ by formulation (21).

**14**     Update network parameter $\theta$ by formulation (22).

**15 end**

**16** Return the network parameter $\theta$.

---

to identify $\hat{\boldsymbol{x}}$ from real samples. Once the discrepancy between reconstructed and real training images is obtained, the parameters of $\mathcal{E}$, $\mathcal{G}_l$ and $\mathcal{G}_g$ can be updated by back-propagation. The semantic encoder and semantic decoder are trained with learning rate $\alpha = 2e^{-4}$ and training epoch $E_p = 70$.

**Stage II: Training stage for semantic bit allocation model:** In this stage, we fix $\mathcal{E}$, $\mathcal{G}_l$, $\mathcal{G}_g$, and $\mathcal{D}$ obtained from stage I, and evoke a randomly initialized RL agent $\pi$ to be trained with reinforcement learning. Specifically, after selecting a quantization level $a^{(m)}$ for the current semantic concept $\mathbf{C}^{(m)}$, the agent will receive a reward $r^{(m+1)}$ indicating whether this action is beneficial. The reward includes a decrease in perceptual loss, the semantic loss, and the rate

**TABLE I:** The classes in the Cityscapes dataset.

| Group | Classes |
|---|---|
| flat | road · sidewalk · parking · rail track |
| human | person · rider |
| vehicle | car · truck · bus · on rails · motorcycle · bicycle · caravan · trailer*+ |
| construction | building · wall · fence · guard rail · bridge · tunnel |
| object | pole · pole group · traffic sign · traffic light |
| nature | vegetation · terrain |
| sky | sky |
| void | ground · dynamic · static |

after taking the action. We train $\pi$ with an off-the-shelf policy gradient algorithm to maximize the cumulated discounted rewards $G$. The procedure for training the RL-based semantic bit allocation model is summarized in Algorithm 2.

**Stage III: Fine tuning stage of the whole model:** In this stage, the semantic encoder $\mathcal{E}$, semantic decoder $\mathcal{G}_g, \mathcal{G}_l$ and RL agent $\pi$ are finetuned together.

## IV. EXPERIMENTAL RESULTS

### A. Simulations Setup

*1) Datasets:* We train and evaluate the RL-ASC model based on the scene parsing and instance segmentation of the Cityscapes dataset [12]. Cityscapes focuses on the semantic understanding of urban street scenes. It is collected from streetscapes in 50 different German cities and consists of 30 classes of objects. There are 2975 images in the training set and 500 images in the validation set; each being annotated with fine semantic labels. We downscale the images and semantic label maps to $256 \times 512$ and conduct testing on the validation set. The classes or the semantic concepts in the Cityscapes dataset are listed as Table I.

*2) Baselines:* We compare the proposed RL-ASC with the engineered codecs BPG [1], JPEG2000 [36], JPEG [47], as well as the deep learning-based codecs DSSLIC [3] and HiFIC [32]. Moreover, we finetuned the DSSLIC with the semantic loss in the semantic segmentation task for fair comparison, and the finetuned model is denoted as DSSLIC-finetuned. The bitrate of the engineered codecs [1], [36], [47] is controlled by the quantization parameters (QP), where

a larger QP means a higher compression ratio. We adopt the officially released version of these codecs and evaluate the performance in the range of bitrates 0 to 0.5 bpp. Particularly, BPG [1] is the current state-of-the-art engineered image compression codec in terms of PSNR. DSSLIC [3] is a layered image compression, where the semantic label of the input image is encoded as the base layer of the bitstream, and the compact representation as well as the residual are encoded as the enhancement layer. The compression ratio of DSSLIC and DSSLIC-finetuned is adjusted by the QP of the enhancement layer. HiFIC [32] combines the GANs with learned compression to achieve high fidelity generative lossy compression, and thus is able to obtain visually pleasing reconstructions that are perceptually similar to the input and operate in a broad range of bitrates. We evaluate the performance of [32] by leveraging the pre-trained models at low bitrate (0.18bpp) and medium bitrate (0.33bpp).

*3) Evaluation Metrics:* To measure the efficiency of the proposed RL-ASC, we evaluate the reconstruction performance from the semantic and perceptual perspectives. The most widely used quality metrics PSNR and SSIM are simple, shallow functions, and fail to account for many nuances of human perception. On the one hand, we evaluate the semantic loss in terms of the performance of downstream tasks such as object detection and semantic segmentation. We adopt mIoU as the objective metric and also illustrate the results of semantic segmentation and object detection as the subjective evaluation. The reconstructed image that suffers less loss in semantic information could yield higer mIoU value, and the maximum value of mIoU is 1.

On the other hand, the metrics Fréchet Inception distance score (FID) [20], and Kernal-Inception distance (KID) [5] (lower better) are consistent with human perception, and thus are employed to evaluate the distance of the reconstructed image and input image in deep feature space. Specifically, KID [5] and FID [20] measure the distribution divergence of the reconstructed images compared with real samples via the Inception network and are widely used to assess sample quality and diversity in the context of GANs. Moreover, the perceptual performance can also be evaluated subjectively by user study, and a better compression approach can yield real, natural, and visual pleasant reconstructions even at a low bitrate.

*4) Compression Modes:* We train three bitrate models of the proposed RL-ASC: the low bitrate (0.08 bpp), the medium bitrate (0.16 bpp), and the high bitrate (0.32 bpp). The bitrate is adjusted by the channels dimension $n$ of the feature map $\boldsymbol{f}$ that is produced by the feature

**TABLE II:** BD-mIoU and BD-rate relative to the baselines for semantic segmentation tasks.

| Metric | DSSLIC [3] | DSSLIC-finetune | HiFiC [32] | BPG [1] | J2K [36] | JPEG [47] | Simplified RL-ASC |
|---|---|---|---|---|---|---|---|
| BD-mIoU | 0.278 | 0.145 | -0.004 | 0.227 | 0.343 | 0.452 | 0.0261 |
| BD-rate | -97.713% | -60.769% | 9.951% | -89.682% | -99.688% | -100% | -31.361% |

extraction network $\mathcal{E}$. We set $n = 16, 32, 64$, $w = W/8$ and $h = H/8$, and therefore the dimensions of $\boldsymbol{f}$ correspond to the three modes are $16 \times W/8 \times H/8$, $32 \times W/8 \times H/8$, and $64 \times W/8 \times H/8$, respectively. Also, the dimension of the downscaled semantic mask $\boldsymbol{s}_d^{(m)}$ is $W/8 \times H/8$. Additionally, we set $Q = 6$, so that the RL agent $\pi$ can choose six quantization levels for different semantic concepts, and higher quantization level means small distortion.

To evaluate the effect of the RL-based semantic bit allocation model, we conduct the ablation study by removing the RL agent $\pi$ from the proposed RL-ASC model. Such an ISC system lacks the adaptive bit allocation ability and thus is denoted as simplified RL-ASC. Specifically, simplified RL-ASC encodes each semantic concept with equal precision by the highest quantization level. The three bitrate models for simplified RL-ASC are low bitrate (0.11 bpp), medium bitrate (0.22 bpp), and high bitrate (0.44 bpp), respectively.

*B. Semantic Performance*

We compare the semantic performance of the proposed RL-ASC and the simplified RL-ASC with the baseline codecs at different bitrates. To validate the effectiveness of the proposed RL-ASC, we apply the proposed method in two downstream semantic tasks: semantic segmentation and object detection.

*1) Objective Quality:* The pretrained PSPNet [53] is adopted as the semantic segmentation model to obtain the semantic label $\hat{\boldsymbol{s}}$ of the reconstructed image. The semantic fidelity can be measured by the consistency between $\hat{\boldsymbol{s}}$ and the ground truth $\boldsymbol{s}$ in terms of mIoU. The mIoU performance of different image codecs in different bitrates is shown in Fig. 5. Cityscapes dataset only contains the ground truth label for semantic segmentation and lacks the ground truth label for other intelligent tasks such as object detection and image classification. Therefore, we cannot obtain the mAP metric of object detection task for objective measurement.
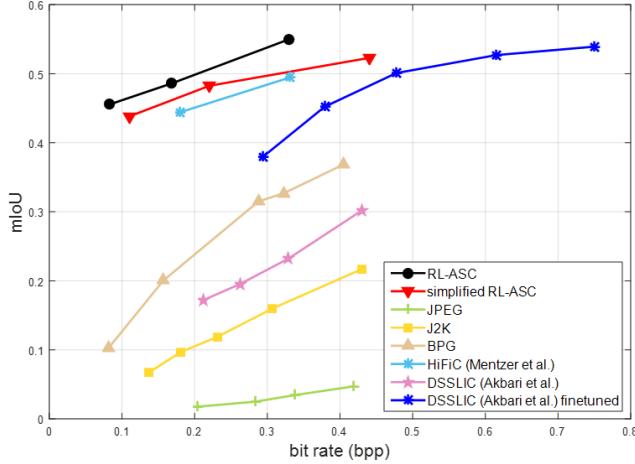
**Fig. 5:** The semantic performance in terms of mIoU of different image codecs on semantic segmentation task. The higher value means better performance.

According to Fig. 5, the semantic fidelity increases with higher bitrates. The proposed RL-ASC outperforms JPEG [47], J2K [36], BPG [1], DSSLIC [3] and DSSLIC-finetuned by a large margin in terms of mIoU, which validates the effectiveness of the proposed task-driven coding manner. The baseline HiFiC [32] approximates the performance of the proposed RL-ASC, while it fails to achieve an extremely low bitrate ($< 0.1$ bpp). The RL-ASC achieves higher mIoU performance compared to the simplified RL-ASC. Note that the DSSLIC finetuned by the semantic loss boosts the performance on the intelligent task in terms of mIoU.

We utilize the Bjontegaard metric [6] to evaluate the coding efficiency of the proposed RL-ASC concerning the baselines. Inspired from [28], we propose BD-mIoU to consider the relative differences between two codecs under equal bitrate in task-related accuracy. BD-mIoU calculates the average mIoU difference between two rate-semantic curves over an interval and the BD-rate represents the average bitrate reduction under the equivalent task-related accuracy. As shown in TABLE. II, the proposed RL-ASC method can achieve the same mIoU with more than 60% bitrate savings on average compared with the deep learning based methods [3], [32] and DSSLIC-finetuned. Under the same bit cost, the proposed RL-ASC method can remarkably improve the mIoU performance compared to the baselines and simplified RL-ASC.
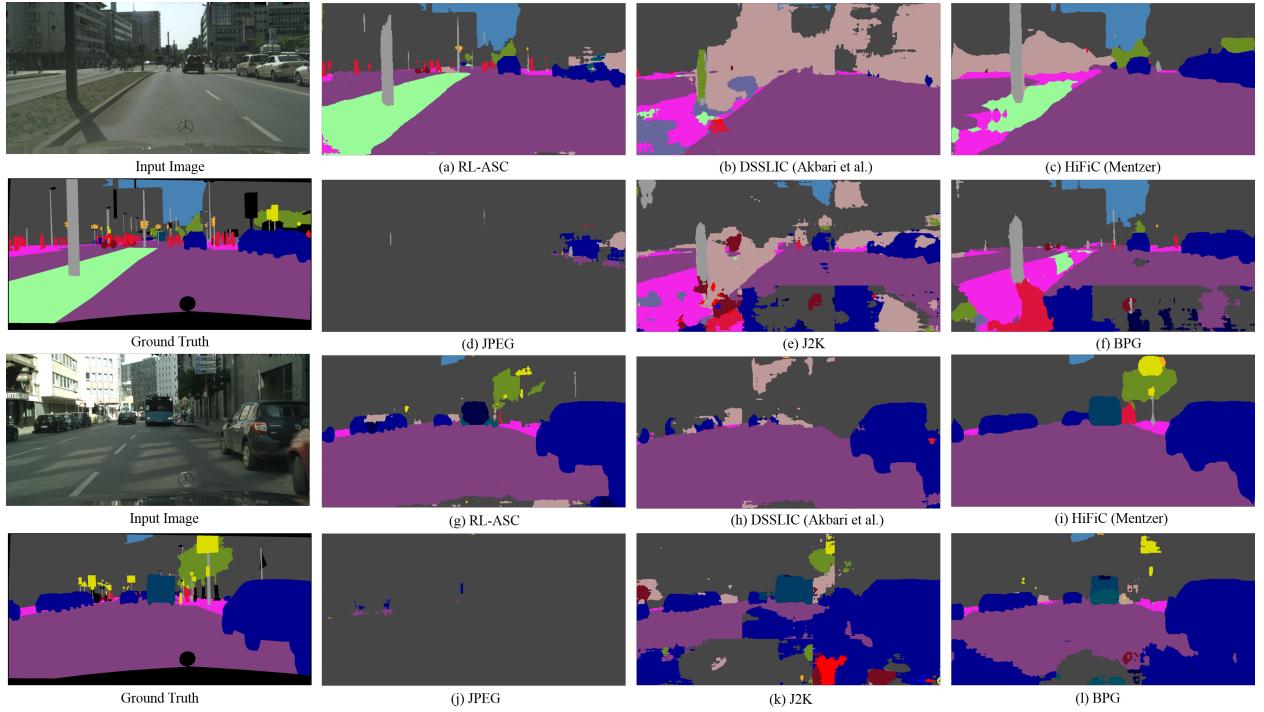
**Fig. 6:** Examples of image coding on downstream semantic segmentation. The first column contains the two randomly selected input image before coding and the corresponding ground truth labels. The rest columns are the semantic label maps of decompressed images of different codecs under similar bitrate (0.33 bpp).

*2) Subjective Quality:* In this section, we visualize the semantic segmentation and object detection results of the decompressed images to validate the semantic fidelity subjectively. To ensure a fair comparison, the proposed RL-ASC encodes the image at a bitrate 0.33 bpp while the baselines encode the image at a bitrate equal to or higher than 0.33 bpp. Particularly, the MaskRCNN [18] pretrained on COCO dataset is adopted to detect objects on the reconstructed image.

The input images, ground truths and the semantic label maps predicted on the decompressed images are shown in Fig. 6. We can observe that the semantic concepts of the proposed RL-ASC can be well recognized and localized on the decompressed images (a) and (g). The segmentation results are similar to the ground truths, which validates that the proposed RL-ASC could maintain the overall meanings of the input image and suffer from little semantic information loss. The deep learning-based codecs DSSLIC [3] and HiFiC [32] perform better than the classic engineered codecs [1], [36], [47]. The semantic labels (b), (c), (h), and (i) represent major concepts such as
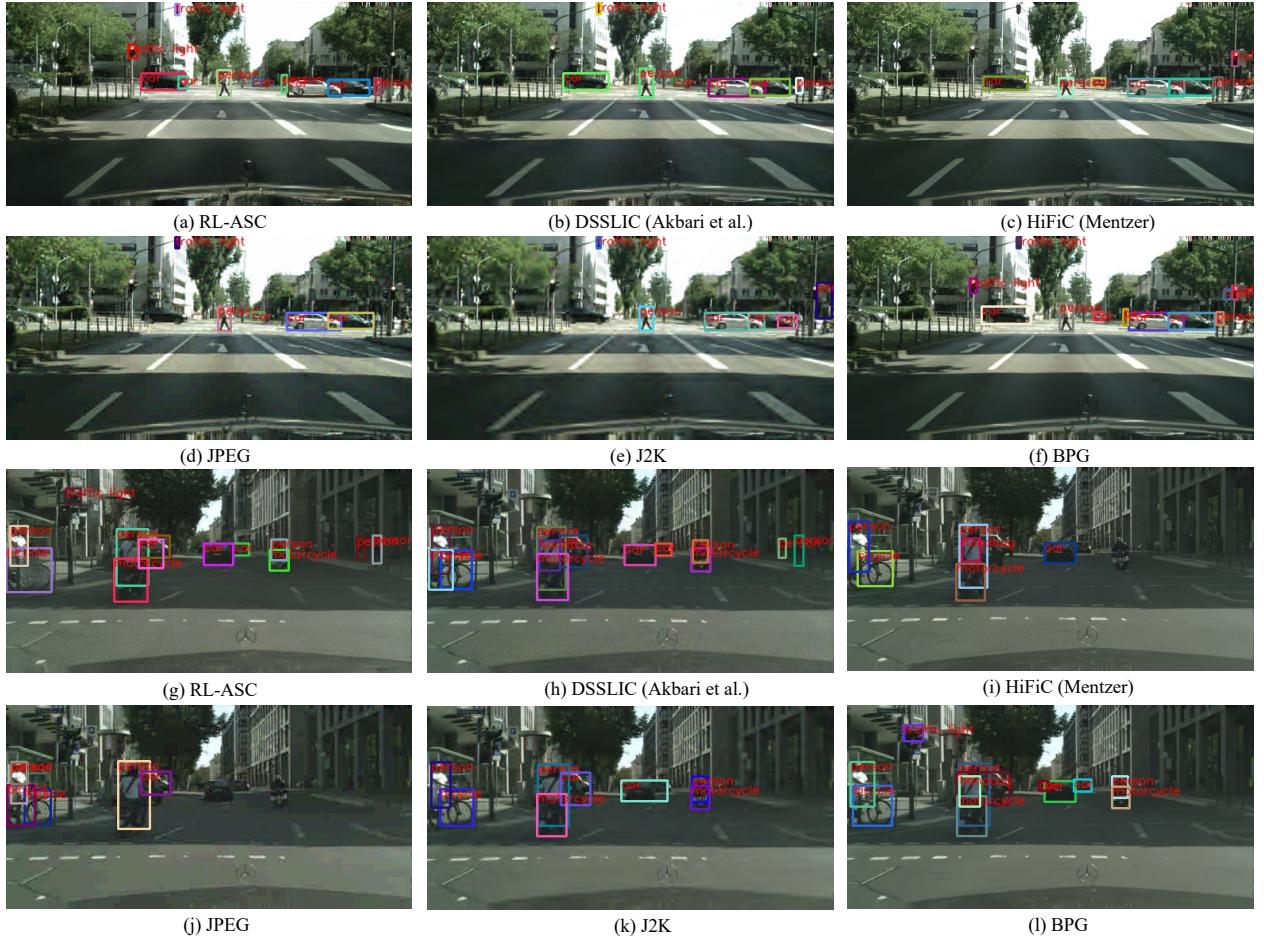
(a) RL-ASC     (b) DSSLIC (Akbari et al.)     (c) HiFiC (Mentzer)

(d) JPEG     (e) J2K     (f) BPG

(g) RL-ASC     (h) DSSLIC (Akbari et al.)     (i) HiFiC (Mentzer)

(j) JPEG     (k) J2K     (l) BPG

**Fig. 7:** Examples of image coding on downstream object detection task. Two randomly selected images are given as examples. The bounding boxes and labels of the decompressed images of different image codecs are produced by the pretrained MaskRCNN.

car, building, and road, while the less important information is missing. As shown in (d), (e), (f), (j), (k), and (i), the outputs of JPEG [47], J2K [36], and BPG [1] fail to conduct semantic segmentation task, and the semantic concepts are misinterpreted by the downstream task. It is because the decompressed image of the engineered codecs degrade heavily at low bitrate, suffering from blocky, blurring or ringing artifact.

The object detection performance on the decompressed images of the proposed RL-ASC as well as the baselines are illustrated in Fig. 7. It is observed that the decompressed images (a) and (g) of the proposed RL-ASC preserve the objects comprehensively. Small objects such as traffic lights and overlapped cars can be detected properly on (a) and (g), which validates the incredible
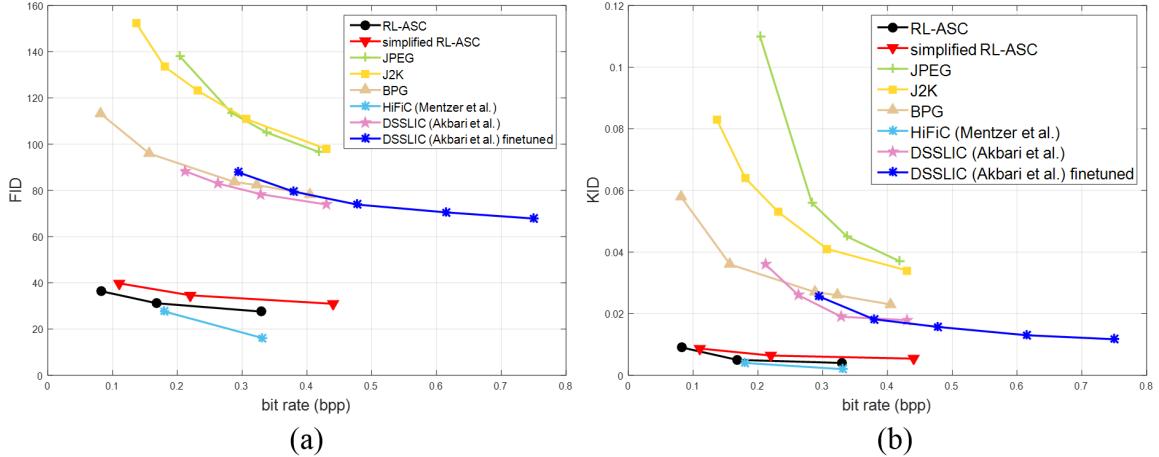
**Fig. 8:** The perceptual performance of different image codecs in terms of (a) FID and (b) KID. The lower value means better performance.

semantic exchange ability of the proposed RL-ASC method. However, since the baselines attempt to recover the exact pixels while ignoring the global underlying meanings, some reconstructed objects fail to be detected accurately. For instance, the traffic light is not detected on HiFiC decompressed image (c) and DSSLIC decompressed image (h), and larger objects such as cars and persons may even be miss detected on classic engineered codecs [1], [36], [47].

### C. Perceptual Performance

We compare the perceptual performance of the proposed RL-ASC and simplified RL-ASC with the baselines at different bitrates.

*1) Objective Quality:* The perceptual performance of different image codecs in terms of FID [20] and KID [5] under different bitrates is illustrated in Fig. 8. The proposed RL-ASC is comparable with HiFiC [32] and outperforms other baselines by a large margin. This can be interpreted as the proposed method incorporates GANs architecture and adopts adversarial loss that enforces the reconstructed image to be natural and realistic, which also validates the effectiveness of the well-designed semantic encoder and semantic decoder. In addition, the RL-based semantic bit allocation model results in convincing increase on perceptual performance under the same bitrate compared to the simplified RL-ASC. There is a moderate degradation for the finetuned DSSLIC compared to the original DSSLIC model in terms of FID and KID. This
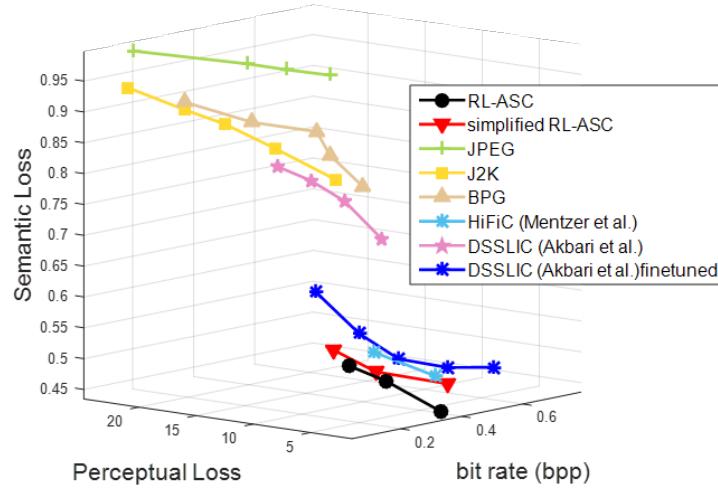
**Fig. 9:** The rate-semantic-perceptual curves of different image codecs. At a certain bitrate, lower semantic loss and lower perceptual loss means better performance.

can be interpreted that the finetuned model emphasizes more on semantic fidelity and sacrifice perceptual performance.

We can further illustrate the triple trade-off rate-semantic-perceptual of the proposed RL-ASC method and the baselines in Fig. 9. In particular, the semantic loss is defined as (1-mIoU), considering the semantic segmentation task. The perceptual loss is defined as the weighted addition of FID and KID values, where a lower value means better performance. As shown in Fig. 9, the proposed RL-ASC achieves lower semantic and perceptual loss compared to the baselines at equal bitrate by a large margin, which validates the effectiveness of the proposed method in semantic exchange and visual performance. Note that the deep learning based methods RL-ASC, [3], [32] outperform the classic engineered codecs, which can be interpreted that the deep features account for better image understanding. Moreover, the DSSLIC finetuned by the semantic loss achieves great progress in this triplet loss, compared to the original DSSLIC model.

*2) Subjective Quality:* The reconstructed images of different image codecs, as well as the original randomly selected input image, are shown in Fig. 10. To ensure a fair comparison, the
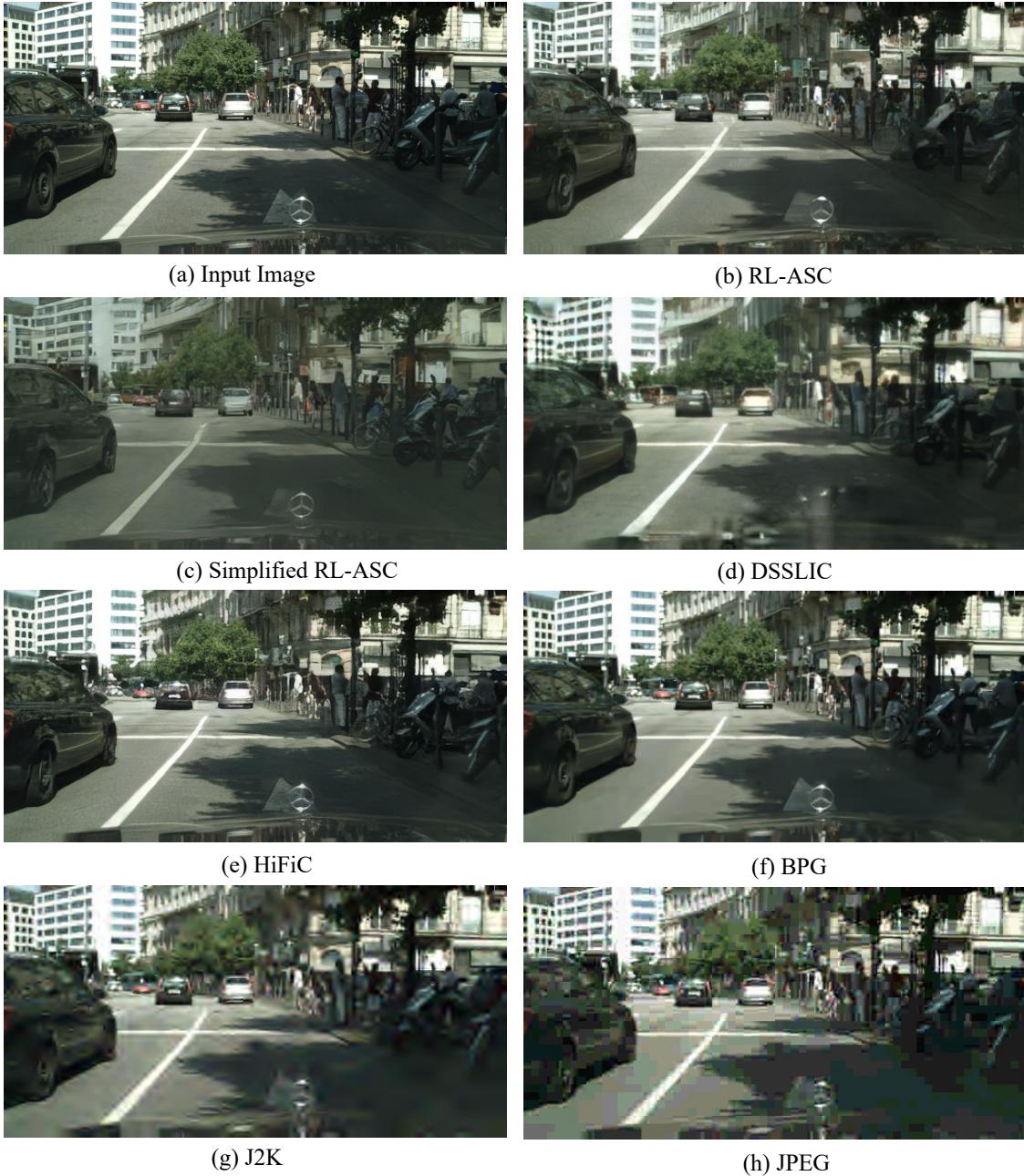
(a) Input Image

(b) RL-ASC

(c) Simplified RL-ASC

(d) DSSLIC

(e) HiFiC

(f) BPG

(g) J2K

(h) JPEG

**Fig. 10:** A randomly selected input image and the reconstructed images of different image codecs at similar bitrate (0.16 bpp).

RL-ASC encodes the image at a bitrate 0.16 bpp while the baselines encode the image at an equal or higher bitrate. It can be observed that the decompressed image (b) of the proposed RL-ASC is almost indistinguishable from the input image (a) even at such a low bitrate. Compared to the simplified RL-ASC (c), the bits in (b) concentrate on salient semantic concepts and therefore
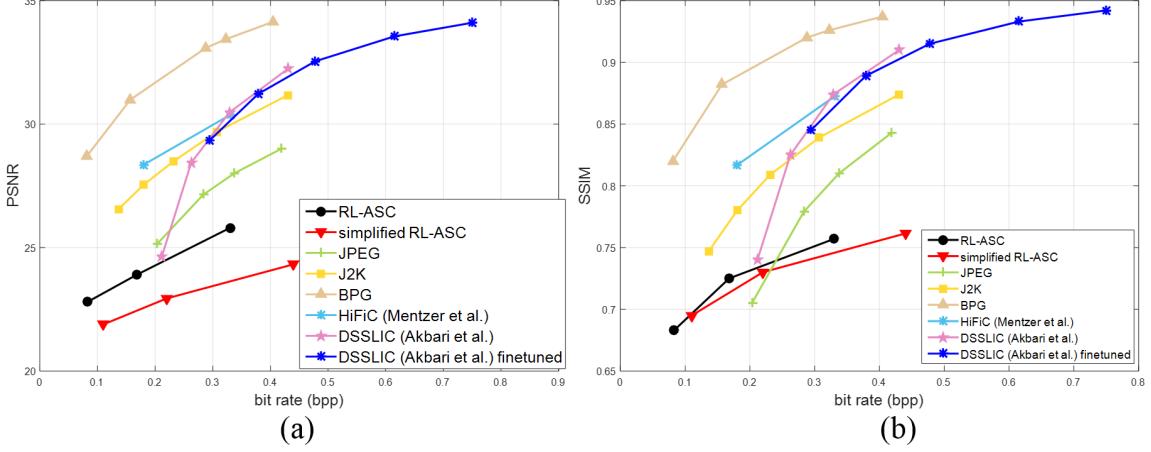
**Fig. 11:** The rate-distortion performance in terms of PSNR and SSIM of different image codecs. Higher value means better performance.
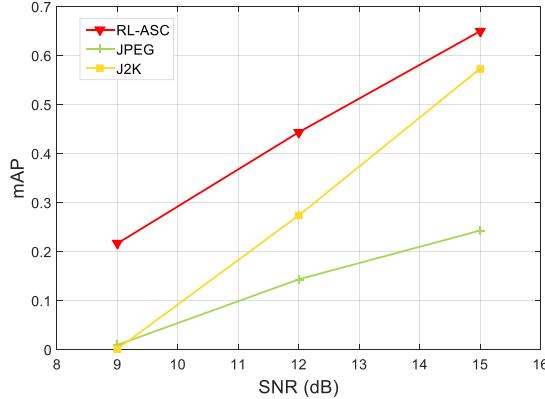


**Fig. 12:** mAP score versus SNR for the same bitrate of compressed images (0.33 bpp) under the AWGN channel.

result in visually pleasant reconstruction. In addition, the baselines suffer from blur, ringing, and blocky artifacts at low bitrate, as shown in (d), (f), (g), and (h).

### D. Rate-Distortion Performance

We evaluate the rate-distortion performance of the proposed RL-ASC, simplified RL-ASC, and the baselines in terms of PSNR and SSIM, which are the most widely used metric in the traditional image coding system that attempts to minimize symbol error. As shown in Fig. 11, the baselines achieve better performance since they are optimized with the PSNR or SSIM metric.

On the contrary, the proposed method is tolerable to pixel errors and does not attempt to ensure local consistency. It can be observed that the RL-ASC outperforms the simplified RL-ASC since the former represents complex objects with higher precision and the rest of the image is simple to be encoded. Moreover, a moderate degradation is occurred in DSSLIC-finetuned compared to DSSLIC in terms of PSNR and SSIM, which validates the trade-off between the classic pixel level loss and the novel semantic loss.

*E. Anti-Noise Performance*

In this last example, we consider the Additive White Gaussian Noise (AWGN) channel and evaluate the robustness of the proposed RL-ASC as well as the classical codecs [36], [47] to physical noise. Fig. 12 shows the object detection performance on the decompressed images in terms of mAP score in different SNRs. Note that the same compression rate (0.33 bpp) is adopted for different methods. The proposed RL-ASC outperforms the baselines by a large margin, which demonstrates the robustness of the proposed method in ISC scenario. The baselines lost most of the information in $SNR = 9$, and the reconstructed image is uninterpretable with mAP $\approx 0$. For $SNR = 15$ or higher, the effect of the physical noise is ignorable, since the mAP accuracy approximates that of the ideal channel condition. In this case, the object detection result obtained by the pretrained MaskRCNN on the input image is deemed as the ground truth label. A more similar detection result on the reconstructed image leads to higher mAP score.

## V. CONCLUSION

In this paper, we considered the ISC system and presented a deep learning-based semantic image coding approach that interprets and encodes images beyond pixel level. We first proposed the novel rate-semantic-perceptual criterion to integrate the semantic fidelity and perceptual quality in the optimization process of the semantic coding. Accordingly, we designed the semantic concept as the novel representation unit and proposed a convolutional semantic encoder to extract semantic information. Driven by the semantic analysis task such as object detection or semantic segmentation, an RL-based semantic bit allocation model is presented to realize the optimization criterion and encode each semantic concept with adaptive quantization. At the receiver side, a generative semantic decoder that adopts attention model to fuse the local and global features

is designed to reconstruct the semantic concepts. With the extracted semantic information, the proposed RL-ASC can facilitate multiple vision tasks in the semantic communication scenario. We compared the decompressed samples of the proposed approach with that of the baselines and showed that FID, KID, and mIoU can be valuable tools to better predict human preferences and the efficiency of semantic exchange. The experiments demonstrated the ability of RL-ASC to produce reconstructions with high semantic similarity, naturalness, and remarkably reduced transmission data amount. Also, the proposed RL-ASC is robust to noise in AWGN channel.

## REFERENCES

[1] https://bellard.org/bpg/.

[2] Eirikur Agustsson, Michael Tschannen, Fabian Mentzer, Radu Timofte, and Luc Van Gool. Generative adversarial networks for extreme learned image compression. In *Int. Conf. Comput. Vis. (ICCV)*, pages 221–231, 2019.

[3] Mohammad Akbari, Jie Liang, and Jingning Han. Dsslic: deep semantic segmentation-based layered image compression. In *2019-2019 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, pages 2042–2046. IEEE, 2019.

[4] Johannes Ballé, Valero Laparra, and Eero P. Simoncelli. End-to-end optimized image compression. In *5th Int. Conf. on Learning Representations, ICLR 2017*, 2017.

[5] Mikołaj Bińkowski, Dougal J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. *arXiv preprint arXiv:1801.01401*, 2018.

[6] Gisle Bjøntegaard. Calculation of average PSNR differences between RD-curves. In *ITU-T VCEG-M33, April,2001*, 2001.

[7] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 6228–6237, 2018.

[8] Eirina Bourtsoulatze, David Burth Kurka, and Deniz Gündüz. Deep joint source-channel coding for wireless image transmission. *IEEE Trans. on Cogn. Commun. Netw.*, 5(3):567–579, 2019.

[9] Zhibo Chen and Tianyu He. Learning based facial image compression with semantic fidelity metric. *Neurocomputing*, 338:16–25, 2019.

[10] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *IEEE Conference Conf. Comput. Vis. and Pattern Recog., (CVPR)*, pages 7936 − 7945, 2020.

[11] Cisco. Cisco visual networking index: Global mobile data traffic forecast update 2017-2022. 2017.

[12] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 3213–3223, 2016.

[13] Samuel Dodge and Lina Karam. Understanding how image quality affects deep neural networks. In *Int. Conf. Qual. Multimed. Exp. (QoMEX)*, pages 1–6, 2016.

[14] Nariman Farsad, Milind Rao, and Andrea Goldsmith. Deep learning for joint source-channel coding of text. In *2018 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2326–2330, 2018.

[15] R. Gong, X. Liu, S. Jiang, T. Li, and J. Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. In *IEEE Int. Conf. on Comput. Vis. (ICCV)*, 2019.

[16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Adv. Neural Inform. Process. Syst.*, 27, 2014.

[17] Klemen Grm, Vitomir Štruc, Anais Artiges, Matthieu Caron, and Hazım K. Ekenel. Strengths and weaknesses of deep learning models for face recognition against image degradations. *IET Biom.*, 7(1):81–89, 2018.

[18] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask R-CNN. In *Proceedings of the IEEE Int. Conf. on Comput. Vis. (ICCV)*, Oct 2017.

[19] T. He, S. Sun, Z. Guo, and Z. Chen. Beyond coding: Detection-driven image compression with semantically structured bit-stream. In *2019 Picture Coding Symposium (PCS)*, 2019.

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Neural Info. Process. Systems (NIPS)*, page 6629–6640, 2017.

[21] Yueyu Hu, Wenhan Yang, and Jiaying Liu. Coarse-to-fine hyper-prior modeling for learned image compression. In *AAAI*, 2020.

[22] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Eur. Conf. Comput. Vis.*, pages 694–711. Springer, 2016.

[23] Nick Johnston, Damien Vincent, David Minnen, Michele Covell, Saurabh Singh, Troy Chinen, Sung Jin Hwang, Joel Shor, and George Toderici. Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks. In *IEEE Conf. Comput. Vis. and Pattern Recog., (CVPR)*, pages 4385 – 4393, 2018.

[24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *Computer Science*, 2014.

[25] David Burth Kurka and Deniz Gündüz. Bandwidth-agile image transmission with deep joint source-channel coding. *IEEE Trans. on Wireless Commun.,*, 2020.

[26] Jooyoung Lee, Seunghyun Cho, and Seung-Kwon Beack. Context-adaptive entropy model for end-to-end optimized image compression. In *7th International Conference on Learning Representations, ICLR*, 2019.

[27] Binglin Li, Linwei Ye, Jie Liang, Yang Wang, and Jingning Han. Region-of-interest and channel attention-based joint optimization of image compression and computer vision. *Neurocomputing*, 500:13–25, 2022.

[28] Xin Li, Jun Shi, and Zhibo Chen. Task-driven semantic coding via reinforcement learning. *IEEE Trans. Image Process.*, 30:6307–6320, 2021.

[29] Dong Liu, Haochen Zhang, and Zhiwei Xiong. On the classification-distortion-perception tradeoff. *Advances in Neural Information Processing Systems*, 2019.

[30] Kun Lu, Rongpeng Li, Xianfu Chen, Zhifeng Zhao, and Honggang Zhang. Reinforcement learning-powered semantic communication via semantic similarity. *arXiv preprint arXiv:2108.12121*, 2021.

[31] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L Van Gool. Conditional probability models for deep image compression. In *IEEE Conf. Comput. Vis. and Pattern Recog. (CVPR)*, 2018.

[32] Fabian Mentzer, George Toderici, Michael Tschannen, and Eirikur Agustsson. High-fidelity generative image compression. *Advances in Neural Information Processing Systems*, 2020.

[33] David Minnen, Johannes Balle, and George. Toderici. Joint autoregressive and hierarchical priors for learned image compression. In *Advances in Neural Information Processing Systems.*, volume 31, 2018.

[34] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 2337–2346, 2019.

[35] Zhijin Qin, Xiaoming Tao, Jianhua Lu, and Geoffrey Ye Li. Semantic communications: Principles and challenges. *CoRR*, abs/2201.01389, 2022.

[36] Majid Rabbani. JPEG2000: Image compression fundamentals, standards and practice. *J Electron Imaging*, 11(2):286, 2002.

[37] Oren Rippel and Lubomir Bourdev. Real-time adaptive image compression. In Doina Precup and Yee Whye Teh, editors, *34th Int. Conf. Mach. Learn. (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pages 2922–2930. PMLR, 06–11 Aug 2017.

[38] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[39] Shibani Santurkar, David Budden, and Nir Shavit. Generative compression. In *2018 Picture Coding Symp. (PCS)*, pages 258–262. IEEE, 2018.

[40] Claude E. Shannon and W. Weaver. The mathematical theory of communication. *Philosophical Review*, 60(3), 1949.

[41] M. Song, J. Choi, and B. Han. Variable-rate deep image compression through spatially-adaptive feature transform. In *Int. Conf. on Comput. Vis.*, 2021.

[42] Simeng Sun, Tianyu He, and Zhibo Chen. Semantic structured image coding framework for multiple intelligent applications. *IEEE Trans. Circuits Syst. Video Technol.*, 31(9):3631–3642, 2021.

[43] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017.

[44] G. Toderici, S. M. O'Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar. Variable rate image compression with recurrent neural networks. *Computer Science*, 2015.

[45] George Toderici, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor, and Michele Covell. Full resolution image compression with recurrent neural networks. In *IEEE Conference Conf. Comput. Vis. and Pattern Recog., (CVPR)*, volume 2017-January, pages 5435 – 5443, 2017.

[46] Robert Torfason, Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Towards image understanding from deep compression without decoding. *6th Int. Conf. on Learning Representations, (ICLR)*, 2018.

[47] Gregory K Wallace. The JPEG still picture compression standard. *IEEE Trans. Consum. Electron.*, 38(1):xviii–xxxiv, 1992.

[48] Zhenzi Weng, Zhijin Qin, and Geoffrey Ye Li. Semantic communications for speech signals. *IEEE Int. Conf. on Commun. (ICC)*, 2020.

[49] Huiqiang Xie and Zhijin Qin. A lite distributed semantic communication system for internet of things. *IEEE J. Sel. Areas Commun.*, 39(1):142–153, 2020.

[50] Huiqiang Xie, Zhijin Qin, Geoffrey Ye Li, and Biing-Hwang Juang. Deep learning enabled semantic communication systems. *IEEE Trans. Signal Process.*, 69:2663–2675, 2021.

[51] Jialong Xu, Bo Ai, Wei Chen, Ang Yang, Peng Sun, and Miguel Rodrigues. Wireless image transmission using deep source channel coding with attention modules. *IEEE Trans. Circuits Syst. Video Technol.*, 2021.

[52] Jiahui Yu, Yuning Jiang, Zhangyang Wang, Zhimin Cao, and Thomas Huang. Unitbox: An advanced object detection network. In *Proceedings of the 24th ACM int. conf. on Multimedia*. ACM, 2016.

[53] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 2881–2890, 2017.