

# CADC: Content Adaptive Diffusion-Based Generative Image Compression

## Supplementary Material

### 8. Implementation Details

#### 8.1. Network Structure

The detailed architectures of the main modules are illustrated in Fig. 8, including the main analysis transform  $g_a$ , main synthesis transform  $g_s$ , hyper analysis transform  $h_a$ , hyper synthesis transform  $h_s$ , auxiliary decoder  $g_{aux}$ , uncertainty estimation network  $f_u$ , and the context model. For efficient network construction, we primarily rely on modified versions of InceptionNeXt [62], GatedCNN [61], and their combination (StarBlock).

We also illustrate the overall entropy modeling process in Fig. 9. Our entropy model integrates a hyperprior module with an autoregressive context model. The hyperprior representation  $\mathbf{c}_h$  is derived as:

$$\mathbf{z} = h_a(\mathbf{y}), \quad \hat{\mathbf{z}} = \lfloor \mathbf{z} \rfloor, \quad \mathbf{c}_h = h_s(\hat{\mathbf{z}}), \quad (14)$$

where  $\lfloor \cdot \rfloor$  denotes rounding to the nearest integer. Here,  $\mathbf{y}$  has 320 channels with a spatial downsampling factor of 64, while  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  also have 320 channels but with a higher spatial compression ratio of  $256\times$ . To balance coding performance and computational efficiency, we employ a 4-step autoregressive process based on quadtree partitioning [35]. As shown in Fig. 9, this process estimates the Gaussian parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$  for the quantized latent  $\hat{\mathbf{y}}$ . Subsequently, arithmetic coding is applied to encode or decode  $\hat{\mathbf{y}}$  into/from the bitstream.

For BLIP, we use the lightweight *blip-image-captioning-base* model (247.41M parameters, with inputs resized to  $384 \times 384$  by default).

To enable a fair evaluation of our three proposed methods, we construct a variant StableCodec by removing its original VAE encoder and auxiliary encoder, and replacing the main analysis/synthesis transforms, hyper analysis/synthesis transforms, and context model with our implementations. This reduces computational complexity and eliminates architectural discrepancies that could otherwise confound the assessment of our contributions. We refer to this variant as our baseline  $M_0$ .

#### 8.2. Training Process

We train our codec based on the standard rate-distortion loss:

$$\mathcal{L} = \lambda \mathcal{R} + \mathcal{D}, \quad (15)$$

where  $\mathcal{R}$  denotes the bitrate,  $\mathcal{D}$  is the distortion measure, and  $\lambda$  is a Lagrange multiplier that balances the two terms. We employ a two-stage training strategy [66]. In the first stage, a base model is trained using a relatively small  $\lambda_{\text{base}}$ .

In the second stage, this pre-trained model is fine-tuned with a larger  $\lambda_{\text{target}}$  to achieve ultra-low target bitrates. The distortion term  $\mathcal{D}$  incorporates multiple components: MSE, LPIPS (computed with VGG features), a CLIP-based loss  $\mathcal{L}_{\text{CLIP}}$  [66]—defined as the  $\ell_2$ -distance between CLIP embeddings of  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ , an adversarial loss  $\mathcal{L}_{\text{adv}}$ , and our proposed auxiliary reconstruction loss  $\mathcal{L}_{\text{aux}}$ . We use DINOv2 [46] with registers [14] as the discriminator backbone [31]. Note that  $\mathcal{L}_{\text{adv}}$  and  $\mathcal{L}_{\text{aux}}$  are only activated in the second training stage. The complete objective is formulated as follows:

$$\text{Stage I : } \arg \min_{\theta} L_1 = \lambda_{\text{base}} \mathcal{R} + \mathcal{D}_1,$$

$$\text{Stage II : } \arg \min_{\theta} L_2 = \lambda_{\text{target}} \mathcal{R} + \mathcal{D}_2,$$

$$\begin{aligned} \mathcal{D}_1 &= d_1 \text{MSE}(x, \hat{x}) + d_2 \text{LPIPS}(x, \hat{x}) + d_3 \mathcal{L}_{\text{CLIP}}(x, \hat{x}) \\ \mathcal{D}_2 &= d_1 \text{MSE}(x, \hat{x}) + d_2 \text{LPIPS}(x, \hat{x}) + d_3 \mathcal{L}_{\text{CLIP}}(x, \hat{x}) \\ &\quad + d_4 \mathcal{L}_{\text{aux}}(x, \hat{x}_{\text{aux}}) + d_5 \mathcal{L}_{\text{adv}}, \end{aligned} \quad (16)$$

where  $\theta$  denotes all trainable parameters in the codec, and  $d_1$ – $d_5$  are weighting coefficients that balance the distortion terms.

### 9. More Experimental Results

#### 9.1. More Quantitative Results

A comprehensive evaluation of rate-distortion performance is conducted on Kodak, the validation set of DIV2K [2], and the test set of CLIC 2020 Professional [57], comparing various codecs (HiFiC [43], DiffEIC [38], GLC [26], DLF [59], ResULIC [28], MKIC [18], OSCAR [21], StableCodec [66]) across six metrics: LPIPS, DISTs, FID, KID, MS-SSIM, and PSNR. Figure 10 shows that our codec consistently delivers state-of-the-art perceptual quality, particularly under ultra-low bitrate conditions.

#### 9.2. More Qualitative Results

Following the quantitative analysis, we also provide a more detailed subjective comparison of the three top-performing codecs: our codec, StableCodec [66], and DLF [59], under ultra-low bitrate conditions. Visual comparisons are presented in Figs. 11 to 14. The results show that our codec consistently achieves the best visual quality, corroborating the quantitative findings and confirming its superior perceptual performance at ultra-low bitrates.

#### 9.3. Runtime Comparison

Table 2 compares the runtime efficiency of our codec against the variant StableCodec and DLF on the Kodak



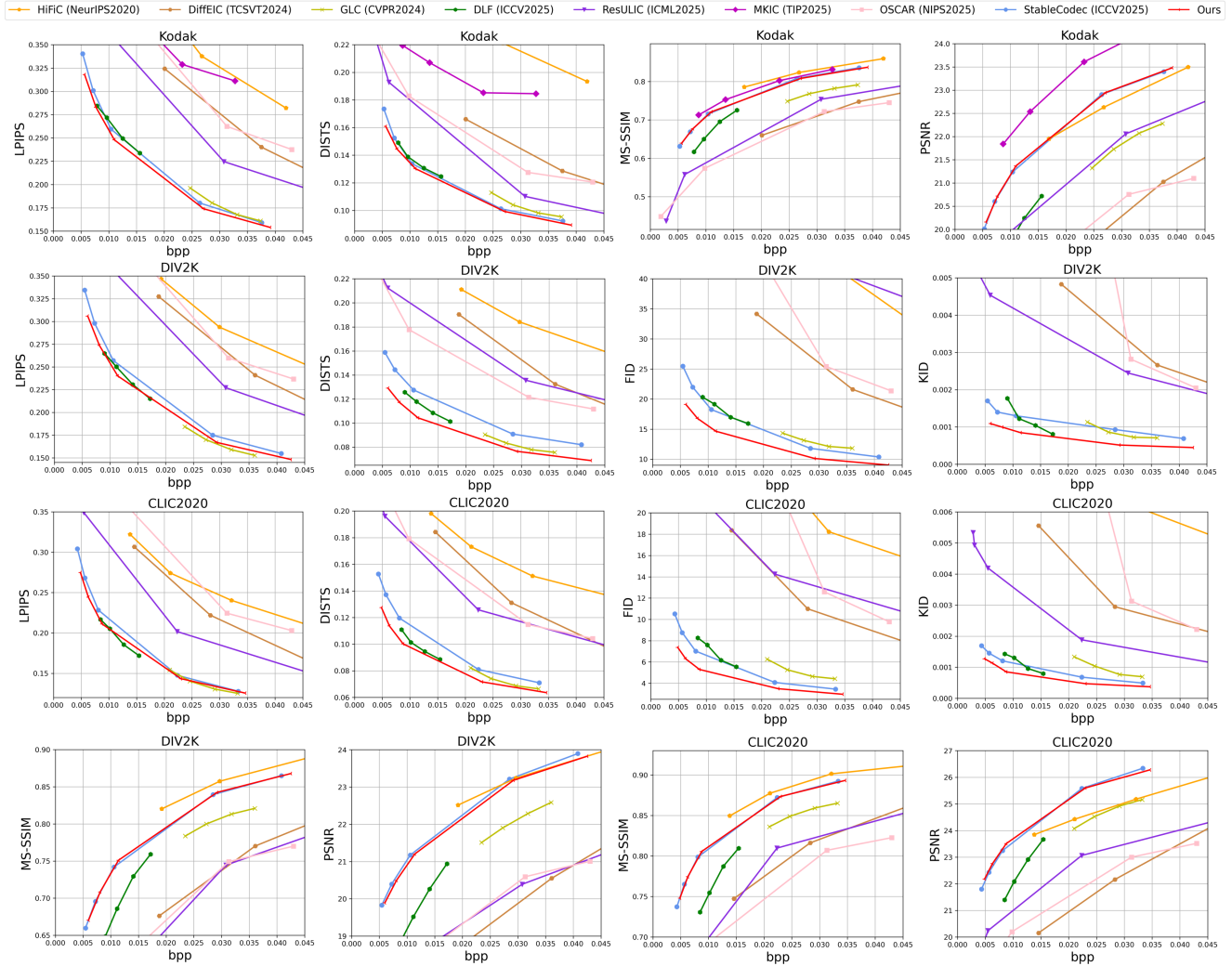


Figure 10. All rate-distortion curves of different generative image codecs on Kodak, the validation set of DIV2K, and the test set of CLIC 2020 Professional in terms of LPIPS, DISTS, FID, KID, MS-SSIM, and PSNR metrics.

#### 9.4. User Study

To comprehensively evaluate the perceptual quality of reconstructed images at ultra-low bitrates, we conduct a user study on the Kodak dataset, comparing our codec with two top-performing codecs: StableCodec [66] and DLF [59]. The study follows a top-1 preference protocol, in which participants are asked to select the reconstruction they perceive as most visually consistent with the ground-truth image. Each participant evaluates 24 randomly ordered test cases. For each case, the ground-truth image is displayed together with the three reconstructed images in a single row of four images, with the order of the methods randomized across trials. Participants are instructed to select the reconstruction that is the most “consistent” with the original image. A total of 25 participants took part in the study, collectively contributing 600 evaluation cases. The results, summarized

Table 3. Top-1 user preference on the Kodak dataset.

Method	Ours	StableCodec	DLF
Bitrate (bpp)	0.0076	0.0072	0.0079
Top-1 Percentage (%)	58.5	29.0	12.5

in Tab. 3, indicate that reconstructions from our codec were preferred in 58.5% of cases. This strong preference underscores its superior perceptual quality as judged by human observers.



Figure 11. Visual examples and comparisons on 2K-resolution images from the test set of CLIC 2020 Professional.

Original

Ours 0.003bpp



StableCodec 0.003bpp

DLF 0.007bpp



Figure 12. Visual examples and comparisons on 2K-resolution images from the test set of CLIC 2020 Professional.

Original

Ours 0.003bpp



StableCodec 0.003bpp

DLF 0.007bpp



Figure 13. Visual examples and comparisons on 2K-resolution images from the validation set of DIV2K.

Original



Ours 0.008bpp



StableCodec 0.008bpp



DLF 0.008bpp



Figure 14. Visual examples and comparisons on 2K-resolution images from the validation set of DIV2K.