

A Novel Framework for Image-to-image Translation and Image Compression

Fei Yang^{a,b,*}, Yaxing Wang^b, Luis Herranz^b, Yongmei Cheng^a and Mikhail G. Mozerov^b

^aSchool of Automation, Northwestern Polytechnical University, Xi'an, 710129, China

^bComputer Vision Center, Universitat Autònoma de Barcelona, Barcelona, 08192, Spain

ARTICLE INFO

Keywords:

Image-to-image translation
Autoencoder
Image compression
Communication

ABSTRACT

Data-driven paradigms using machine learning are becoming ubiquitous in image processing and communications. In particular, image-to-image (I2I) translation is a generic and widely used approach to image processing problems, such as image synthesis, style transfer, and image restoration. At the same time, neural image compression has emerged as a data-driven alternative to traditional coding approaches in visual communications. In this paper, we study the combination of these two paradigms into a joint I2I compression and translation framework, focusing on multi-domain image synthesis. We first propose distributed I2I translation by integrating quantization and entropy coding into an I2I translation framework (i.e. I2Icodec). In practice, the image compression functionality (i.e. autoencoding) is also desirable, requiring to deploy alongside I2Icodec a regular image codec. Thus, we further propose a unified framework that allows both translation and autoencoding capabilities in a single codec. Adaptive residual blocks conditioned on the translation/compression mode provide flexible adaptation to the desired functionality. The experiments show promising results in both I2I translation and image compression using a single model.

1. Introduction

Modern computer vision and image processing heavily rely on deep neural networks and machine learning. One prominent example is image-to-image (I2I) translation, which addresses the problem of learning to transform images from a source domain to a target domain. This general approach has numerous applications in image restoration and enhancement (e.g. colorization, super-resolution, deblurring), but also more complex data-driven transformations where the transformation is learned from data (e.g. style transfer, face attribute modification, scene synthesis, zebra-to-horse translation). Recently, deep neural networks have been also applied to image coding, resulting in the alternative coding paradigm of neural image compression (NIC). These approaches can compete and often surpass the rate-distortion performance of traditional transform coding approaches (e.g. BPG [1]). Image and video coding technology has also significant implications in visual communications and the storage and distribution of video content.

Building upon those two aforementioned research areas, in this paper we study the problem of *distributed I2I translation*, where encoding is performed at the sender side, decoding at the receiver side and the coded representation is either transmitted through a digital communications channel or stored. Thus, in addition to addressing the translation problem, we also aim at obtaining compact binary representations (i.e. bitstreams).

A naive approach to this problem would be translating the image at the sender side and then compressing the result, or translating the reconstructed image at the receiver

side. These approaches have several limitations. First, they require encoding and decoding images twice, once with the translator and once with the image codec, resulting in lower computational efficiency. Similarly, it requires storing two separate codecs (i.e. for translation and autoencoding). Finally, each encoding and decoding pass is a lossy transformation, therefore it is likely that through these four transformations more information is lost, resulting in lower quality in the translation with artifacts, and/or larger bitstreams. In order to address these limitations, we propose the I2Icodec framework (see Figure 1a), which addresses distributed I2I translation with a single encoder and decoder, thus avoiding computational overheads and potential loss of information.

While I2Icodec only requires a single encoder and decoder pair, it cannot perform regular autoencoding (i.e. conventional image coding), which is a desirable functionality in practice. A naive solution is to deploy a regular image codec alongside, but that increases the memory requirements significantly. Thus, to avoid deploying two separate models, we also propose UI2Icodec, a unified framework that can perform both distributed translation and autoencoding in a single model (see Figure 1b).

Note that distributed I2I translation is not limited to *spatially* dislocated encoder and decoder, but it could also be applied to store local files that will be decoded in the future (i.e. *temporally* distributed I2I translation). In this way, a single compact model prevents unnecessary use of computation and memory resources.

In summary, our contributions are: 1) we study the problem of distributed I2I translation, which involves I2I translation under rate constraints; 2) a novel framework for distributed I2I translation (I2Icodec); and 3) a unified framework for distributed I2I translation and autoencoding (UI2Icodec).

*Corresponding author

✉ fyang@cvc.uab.cat (F. Yang)

ORCID(s):

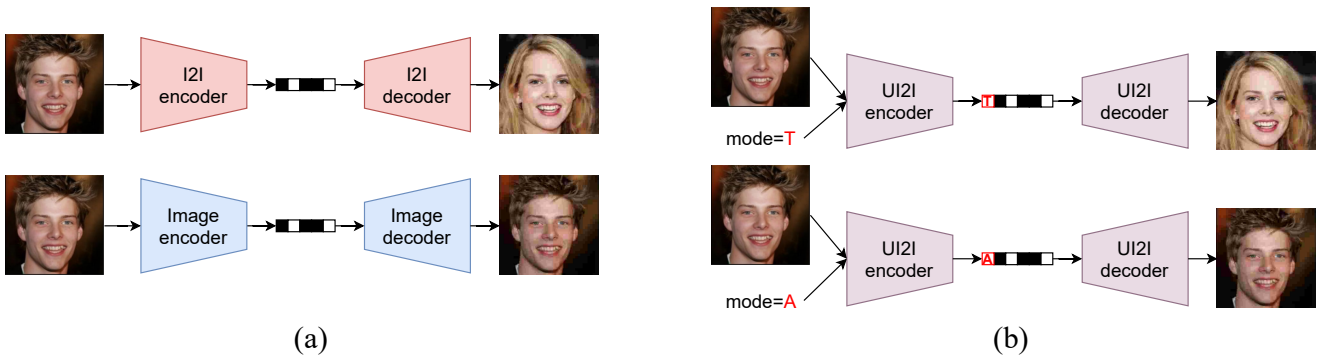


Figure 1: Proposed approaches: (a) distributed I2I translation (*I2Icodec*), alongside a regular image compression codec, and (b) unified translation and autoencoding framework using a single codec (*UI2Icodec*).

2. Related work

2.1. Image-to-image translation

Image-to-image translation has been studied widely in recent years. Paired I2I translation [2, 3, 4, 5, 6] assumes the availability of input-output image pairs. Unpaired I2I translation [7, 8, 9, 10, 11, 12, 13] is a more challenging setting where translations are learned from (unpaired) images from the input and output domain.

Often, a given input image can have multiple plausible translations (e.g. colorization of grayscale images). Multimodal I2I translation (or diverse I2I translation) [14, 15, 16, 17, 18] addresses this problem by disentangling content and style. Style is sampled randomly, ensuring the model generates diverse translations. Early approaches assume only two domains. More recently, multi-domain I2I translation approaches [18, 19, 20] can translate between a range of domains using a single model. In particular, we build upon on StarGAN v2 [19] which provides state-of-the-art translation, including multimodal and multi-domain translation, with content-style disentanglement.

2.2. Neural image compression

Motivated by the success of deep learning, neural image compression [21, 22, 23] has emerged as a new paradigm where codecs are implemented as deep neural networks. Their parameters are directly optimized to minimize a particular combination of rate and distortion over a training dataset, a clear advantage over traditionally engineered transform coding. One key obstacle is quantization and entropy coding, which are non-differentiable operations. In practice, during training, quantization is replaced by a differentiable proxy [22, 23], or soft vector quantization [24], and entropy coding is bypassed, with rate estimated as the entropy of the quantized latent representation. The architecture is based on convolutional autoencoder [22, 23], sometimes with recurrent neural networks [21, 25]. A key component in minimizing the rate is the learnable entropy model. Recent models include hyperpriors [26], autoregressive models [27, 28, 29, 30, 31] and generative models [32, 33]. More recently, variable-rate approaches enable encoding at multiple rate-distortion tradeoffs within

the same model [23, 34, 35]. Aiming at decoding realistic reconstructions even with low rates, some works [36, 37, 38, 39] explore the use of perceptual and adversarial losses during training. Motivated by this, we also use adversarial loss for both translation and autoencoding.

3. Distributed image-to-image translation

We first consider the problem of distributed I2I translation to a target domain. In particular, an image $\mathbf{x} \in \mathcal{X}$ from a source domain label $y_{src} \in \mathcal{Y}_{src}$ is transformed in a compact latent representation \mathbf{z} and encoded into a bitstream \mathbf{b} at the sender side. Following the common practice of disentangling content and style, the expected style of the translated image should also be sent to the receiver side while \mathbf{z} corresponds to the content. Then, the receiver side can decode \mathbf{b} , reconstructing the translated image in a target domain, indicated by the label $y_{tar} \in \mathcal{Y}_{tar}$. The style of the translated image is determined by a style vector \mathbf{s} , either sampled randomly or obtained from a reference image. This disentanglement enables the reconstruction of diverse translations for a given image \mathbf{x} . The style vector \mathbf{s} can be provided by either the sender or the receiver. In the former, the style vector is quantized and included in the bitstream \mathbf{b} and transmitted or stored accordingly (being very compact, the overhead is negligible).

The objective in distributed I2I translation is to obtain successful translations with compact bitstreams.

3.1. I2Icodec framework

Our framework is based on the I2I translation framework of [19], but with the encoder and decoder located separately in the sender and receiver sides, respectively (see Figure 2). In order to achieve this, the framework is augmented with compression capabilities, i.e. quantization and entropy coding. It is composed of *content encoder* E^c , *style encoder* E^s , *mapping network* M , *decoder* G and *discriminator* D^T .

Content encoder The content encoder E^c extracts a latent representation $\mathbf{z} = E^c(\mathbf{x})$ of the content of the image \mathbf{x} . To transmit this representation via a binary channel, the representation \mathbf{z} is quantized (in this paper we use scalar

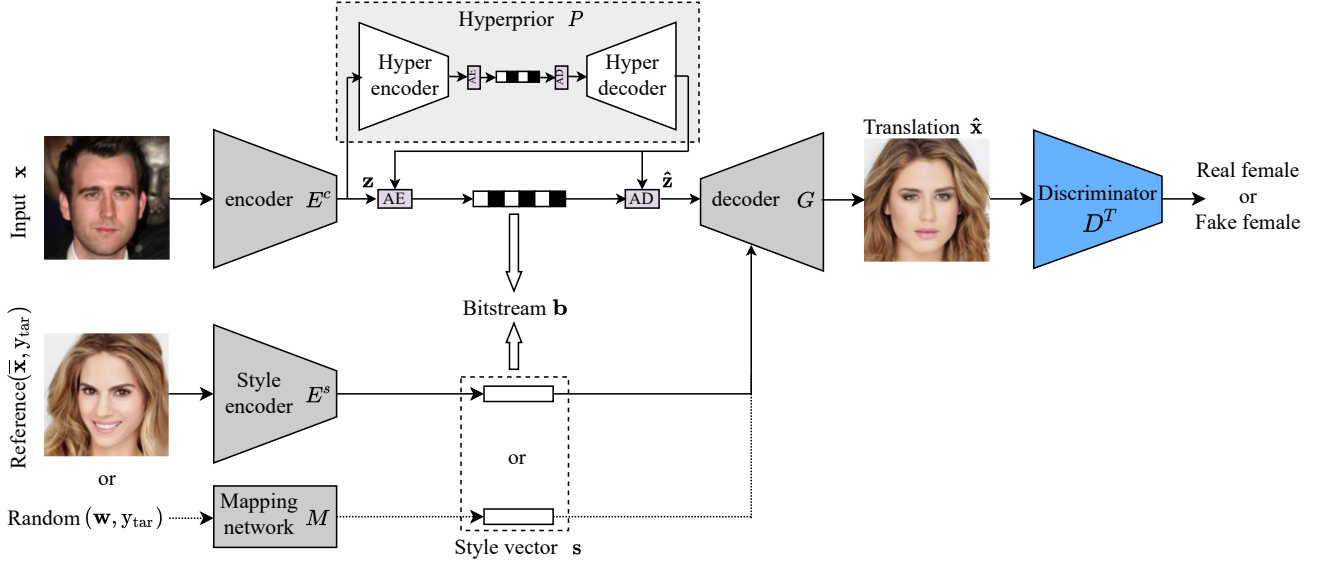


Figure 2: Architecture of the framework for distributed I2I translation (*I2Icodec*). The style vector \mathbf{s} needs to be included in the bitstream when it is provided on the sender side and not when on the receiver side.

quantization) as $\mathbf{q} = Q(\mathbf{z})$ to obtain a discrete-valued representation $\mathbf{q} \in \mathbb{Z}^D$. Then, \mathbf{q} is binarized using entropy coding (arithmetic coding in our case) to reduce statistical redundancy. Quantization is lossy, but entropy coding is not. During training, we replace quantization by uniform noise, and bypass entropy coding, approximating the rate by the entropy of \mathbf{z} .

Mapping network and style encoder. The style \mathbf{s} is used for guiding I2I translation towards a specific style in the target domain. This style code can either be sampled randomly (providing diversity), or obtained from a reference style image. In the former, the mapping network M obtains the domain-specific style representation \mathbf{s} from a domain-independent random style \mathbf{w} as $\mathbf{s} = M(\mathbf{w}, y_{tar})$. Alternatively, the domain-specific style representation can be obtained from a reference image $\bar{\mathbf{x}}$ with the style encoder E^s as $\mathbf{s} = E^s(\bar{\mathbf{x}}, y_{tar})$.

Decoder. The decoder G receives the bitstream and performs entropy decoding and maps back to the real-valued representation $\hat{\mathbf{z}}$. It then generates the reconstructed image from \mathbf{z} and the style \mathbf{s} as $\hat{\mathbf{x}} = G(\hat{\mathbf{z}}, \mathbf{s})$.

Discriminator. Following [19], we use a multi-task discriminator where $D^T(\hat{\mathbf{x}}, y_{tar})$ returns the probability that $\hat{\mathbf{x}}$ is classified in target domain y_{tar} .

Entropy model. We use a learnable hyperprior [26] to model the latent distribution $P(\mathbf{z})$.

3.2. Loss

Our objective during training is to optimize translation while minimizing the rate. Regarding translation, we follow the losses used in [19]. Given an image $\mathbf{x} \in \mathcal{X}$ and its original domain label $y_{src} \in \mathcal{Y}_{src}$, we can obtain its latent representation $\mathbf{z} = E^c(\mathbf{x})$. The I2I translation loss \mathcal{L}_T consists of the following terms.

Adversarial loss. In principle, the adversarial loss forces the translated images to be indistinguishable from real photos. We first generate a random domain-specific target style as $\tilde{\mathbf{s}} = M(\mathbf{w}, \tilde{y}_{tar})$ from a random domain-independent style \mathbf{w} and random target domain $\tilde{y}_{tar} \in \mathcal{Y}_{tar}$. The decoder then synthesizes the translated image as $G(\hat{\mathbf{z}}, \tilde{\mathbf{s}})$. Finally, we employ adversarial loss [40] as

$$\mathcal{L}_{adv} = \mathbb{E}_{\mathbf{x}, y_{src}} [\log D^T(\mathbf{x}, y_{src})] + \mathbb{E}_{\mathbf{x}, \tilde{y}_{tar}, \mathbf{w}} [\log (1 - D^T(G(\hat{\mathbf{z}}, \tilde{\mathbf{s}}), \tilde{y}_{tar}))], \quad (1)$$

where G tries to generate images $G(\hat{\mathbf{z}}, \tilde{\mathbf{s}})$ that look similar to images from domain \tilde{y}_{tar} , while the discriminator D^T aims to distinguish between translated samples $G(\hat{\mathbf{z}}, \tilde{\mathbf{s}})$ and real samples \mathbf{x} from domain y_{src} . G aims to minimize this objective against an adversary D^T that tries to maximize it, i.e., $\min_G \max_{D^T} \mathcal{L}_{adv}(G, D^T, \mathcal{Y}_{src}, \mathcal{Y}_{tar})$.

Style reconstruction. We encourage the decoder to optimize the style representation $\tilde{\mathbf{s}}$ when generating the image $G(\hat{\mathbf{z}}, \tilde{\mathbf{s}})$ with a style reconstruction loss

$$\mathcal{L}_{sty} = \mathbb{E}_{\mathbf{x}, \tilde{y}_{tar}, \mathbf{w}} [\|\tilde{\mathbf{s}} - E^s(G(\hat{\mathbf{z}}, \tilde{\mathbf{s}}), \tilde{y}_{tar})\|_1]. \quad (2)$$

Style diversification. To encourage diversity and prevent mode collapse, we sample and map random pairs of styles $\tilde{\mathbf{s}}_1 = M(\mathbf{w}_1, \tilde{y}_{tar})$ and $\tilde{\mathbf{s}}_2 = M(\mathbf{w}_2, \tilde{y}_{tar})$, using diversity sensitive loss

$$\mathcal{L}_{ds} = \mathbb{E}_{\mathbf{x}, \tilde{y}_{tar}, \mathbf{w}_1, \mathbf{w}_2} [\|G(\hat{\mathbf{z}}, \tilde{\mathbf{s}}_1) - G(\hat{\mathbf{z}}, \tilde{\mathbf{s}}_2)\|_1]. \quad (3)$$

Cycle consistency. To ensure that the domain-invariant structure of the input image \mathbf{x} in the translated image $G(\mathbf{z}, \tilde{\mathbf{s}})$ is preserved we use the cycle consistency mechanism [10]

$$\mathcal{L}_{cyc} = \mathbb{E}_{\mathbf{x}, y_{src}, \tilde{y}_{tar}, \mathbf{w}} [\|\mathbf{x} - G(E^c(G(\hat{\mathbf{z}}, \tilde{\mathbf{s}})), \hat{\mathbf{s}})\|_1], \quad (4)$$

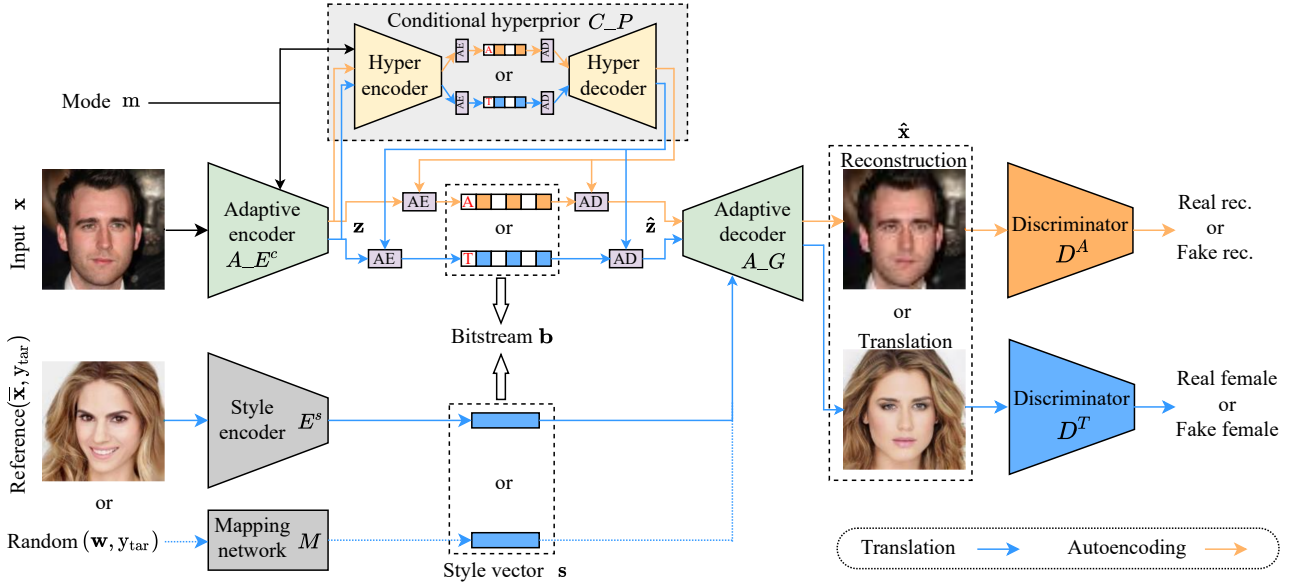


Figure 3: Architecture of the unified framework for regular image compression and I2I translation (UI2Icodec). The mode m signals whether the model runs as autoencoder or as I2I translator (store as 1 bit in the bitstream).

where $\hat{\mathbf{s}} = E^s(\mathbf{x}, y_{src})$ is the style of the input image.

Rate. We estimate the rate as the entropy of the bitstream via modeling the distribution of \mathbf{z} using the entropy model P . This term encourages the model to retain the most important information in a compact representation

$$\mathcal{L}_{rate} = \mathbb{E}_{\mathbf{x}} [-\log(P(\mathbf{z}))] \quad (5)$$

using $\mathbf{z} = E^c(\mathbf{x})$. The final loss is

$$\mathcal{L}_T = \mathcal{L}_{adv} + \gamma_{sty} \mathcal{L}_{sty} - \gamma_{ds} \mathcal{L}_{ds} + \gamma_{cyc} \mathcal{L}_{cyc} + \lambda_T \mathcal{L}_{rate}. \quad (6)$$

4. Unified translation and autoencoding

While the I2Icodec framework can realize distributed I2I translation, being able to recover the original input image (i.e. regular image compression) is equally important in practice. In order to avoid having to deploy two independent codecs (i.e. I2Icodec and autoencoding codec), here we propose a unified framework to transmit an input image and recover either reconstruction of the original image (autoencoding mode) or a translated image (translation mode), in a single model.

4.1. UI2Icodec framework

As shown in Figure 3, we endow the I2Icodec with a switching mechanism controlled via the mode input $m \in \{A, T\}$, which signals the operating mode. In the following, we describe the additional modifications to the I2Icodec framework to implement the joint functionality.

Adaptive encoder and decoder. The content encoder $A_E^c(\mathbf{x}; m)$ and the decoder $A_G(\hat{\mathbf{z}}, \mathbf{s}; m)$ are conditioned on the mode m . When $m = T$, the encoder and decoder operate exactly as the I2Icodec described earlier; When

$m = A$, the model works in autoencoding mode, i.e. normal neural image compression.

Adaptive residual blocks (AdaResBlocks). The switching functionality is implemented by conditioning the residual blocks of the content encoder and decoder on the mode m , via an adaptation unit that modulates intermediate features within the residual block (see Figure 4 (a) and (b)). As shown in Figure 4 (c), adaptation units of the content encoder modulate a given input feature \mathbf{o} as $\mathbf{o}' = \mathbf{u}(m) \odot \mathbf{o} + \mathbf{r}(m)$. The adaptation units in the decoder implement $\mathbf{o}' = \mathbf{u}(\mathbf{s}; m) \odot \mathbf{o} + \mathbf{r}(\mathbf{s}; m)$. In both cases, scale and bias parameters themselves are obtained via linear functions of m (and \mathbf{s} in the decoder).

Conditional entropy model. We condition the hyperprior [26] on the mode, i.e. $C_P(\mathbf{z}; m)$, and one underlying factorized model for translation and another for autoencoding, selected depending on m .

Task-specific discriminators. We use a separate discriminator D^A (see Figure 3) when optimizing the autoencoding task. We found this to be more effective than using a shared discriminator for both tasks.

4.2. Losses

During training we optimize a loss with two terms corresponding to each of the operating modes

$$\mathcal{L} = [m = A] \mathcal{L}_A + [m = T] \mathcal{L}_T, \quad (7)$$

where $[I]$ is the Iverson bracket (1 when I is true, 0 otherwise), \mathcal{L}_T is the loss described in the previous section minimized when $m = T$, and \mathcal{L}_A is the autoencoding loss, minimized for $m = A$. The latter combines the losses of rate-distortion tradeoff and GAN loss as $\mathcal{L}_T = \mathcal{L}_{RD} + \beta \mathcal{L}_{adv2}$, which are introduced in the following.

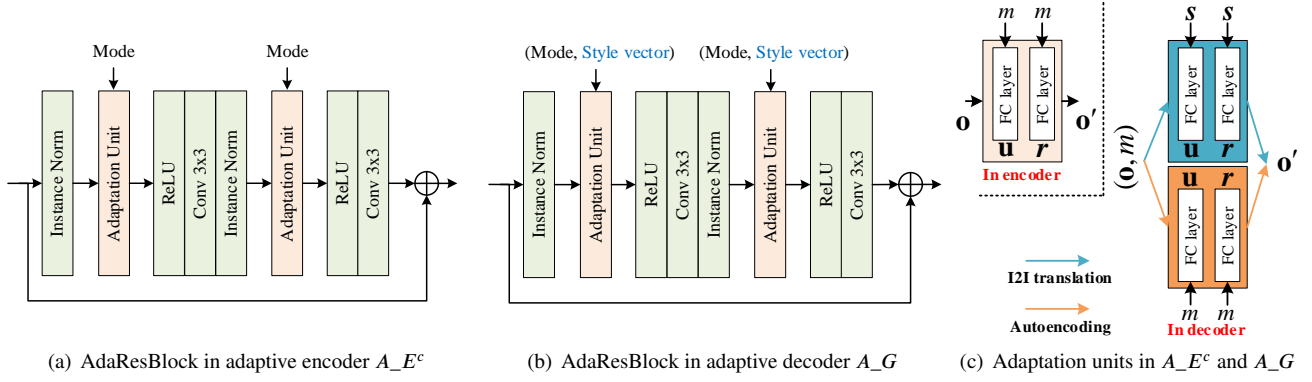


Figure 4: (a-b): Adaptive ResBlock for switching between I2I translation and autoencoding; (c) Details of adaptation units in encoder A_{E^c} and decoder A_G .



Figure 5: Baselines for distributed I2I translation.

Rate-distortion. In autoencoding we optimize a combination of rate and distortion, where the tradeoff is controlled by the parameter λ_A which also decides the final bit-rate of compression model. The loss is

$$\mathcal{L}_{RD} = \mathbb{E}_{\mathbf{x}} [d(\mathbf{x}, \hat{\mathbf{x}}) - \lambda_A \log(C_P(\mathbf{z}; m = A))] \quad (8)$$

where $\mathbf{z} = A_{E^c}(\mathbf{x}; m = A)$ is the latent representation, $\hat{\mathbf{x}} = A_G(\hat{\mathbf{z}}; m = A)$ is the reconstructed image, and $d(\mathbf{x}, \hat{\mathbf{x}})$ is the distortion metric (mean-squared error in our case).

Adversarial loss for generative image compression. We also encourage realism in the reconstructed images by introducing adversarial training with the discriminator D^A , which has been verified that it can help to generate high fidelity images even with an extreme low rate [39]. Similarly to \mathcal{L}_{adv} , we apply the second adversarial loss as

$$\begin{aligned} \mathcal{L}_{adv2} = & \mathbb{E}_{(\mathbf{x}, y_{src})} [\log D^A(\mathbf{x}, y_{src})] \\ & + \mathbb{E}_{(\mathbf{x}, y_{src})} [\log(1 - D^A(\hat{\mathbf{x}}, y_{src}))]. \end{aligned} \quad (9)$$

Then, optimizing $\min_{A_G} \max_{D^A} \mathcal{L}_{adv2}(A_G, D^A, \mathcal{Y}_{src}, \mathcal{Y}_{src})$ will lead A_G with $m = A$ to generate more realistic images.

5. Experiments

In this section, we describe our experimental setup and results. We analyze the effects of the rate constraint of I2Icodec in the translations (see Section 5.1). We also show the results of our unified framework UI2Icodec on both autoencoding and I2I translation (Section 5.2). Ablation study and additional results are shown in Section 5.3.

Datasets. Our experiments are mainly conducted on CelebA-HQ and the animal faces (AFHQ) datasets [19].

As in [19], CelebA-HQ is separated into male and female domains, and AFHQ into cat, dog and wildlife domains. We resized all images to 256×256 for training and comparisons.

Training. For I2Icodec, we train the model minimizing \mathcal{L}_T during 100k iterations. We set $\gamma_{sty} = \gamma_{ds} = \gamma_{cyc} = 1$ which are same with [19], and $\lambda_T \in \{0.01, 0.05, 0.1, 0.3, 0.5\}$ (in Eq. 6) which can lead to different bit-rates for the distributed I2I translation. For UI2Icodec, we first train A_{E^c} and A_G just with distortion loss (mean square error) on autoencoding mode for 50k iterations, then the whole model is jointly trained with another 100k iterations. We use Adam [43] to optimize the model with the loss of Eq. 7 calculated by setting $m = A$ for autoencoding and $m = T$ for I2I translation during each iteration. We set $\lambda_A = [5, 10, 15, 20, 30]$ and $\beta = 1$ to achieve different bit-rates for normal image compression, and fix $\lambda_T = 0.5$ for translation on CelebA-HQ dataset and $\lambda_T = 0.1$ on AFHQ dataset.

Evaluation metrics. We rely on the usual metrics used in I2I translation and image compression. For translation we compute Fréchet inception distance (FID) [44] and learned perceptual image patch similarity (LPIPS) [45] to evaluate quality and diversity, respectively. Autoencoding is evaluated by computing the distortion metrics PSNR, MS-SSIM, and LPIPS. Note that, for translation LPIPS is computed between pairs of translated images, while for autoencoding it is computed between original and reconstructed images. The rate is measured as the bits per pixel (BPP) of the bitstream including the content part output by E^c of I2Icodec or A_{E^c} of UI2Icodec and the style part which is a 64-dimensional vector in Float32 and requires 0.03125 BPP for 256×256 images. Note that the contribution of other parts of

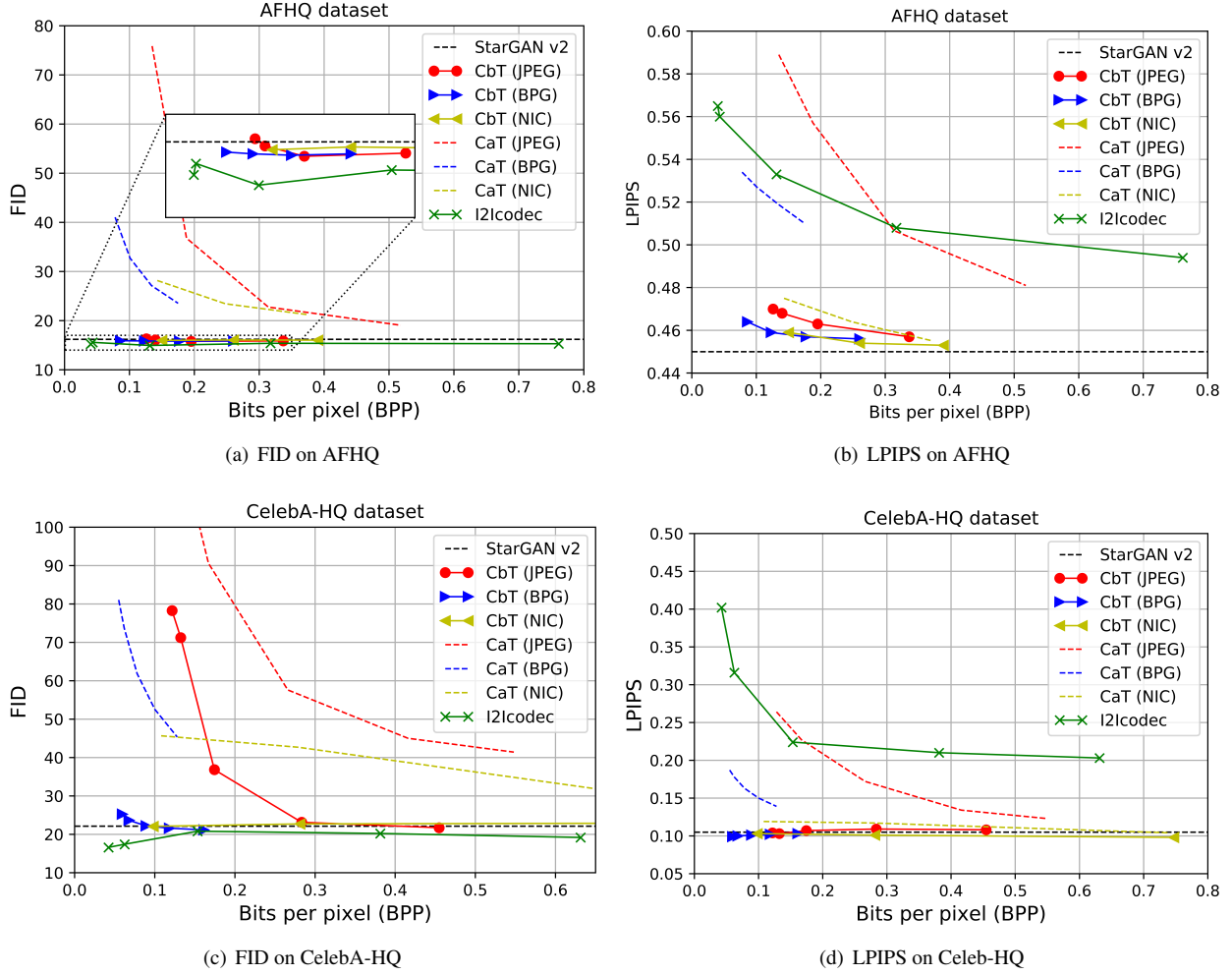


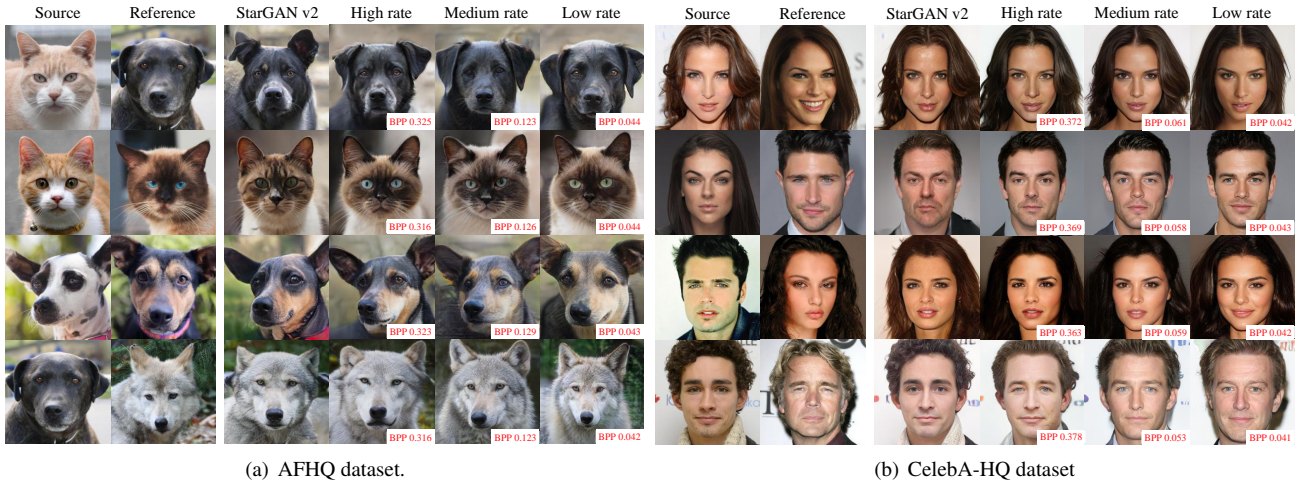
Figure 6: Results of different scheme on AFHQ and CelebA-HQ dataset.

the bitstream (i.e. operation mode) to the rate is negligible and ignored in our results.

5.1. Distributed I2I translation

In this section, we analyze and compare different methods to address distributed I2I translation: (1) compression before translation (CbT): the input image is compressed and then translated after reconstruction, as shown in Figure 5 **left**; (2) compression after translation (CaT): the input image is translated and then compressed with the image codec, as shown in Figure 5 **right**; (3) the proposed I2Icodec (see Figure 2). We use the pre-trained model same with [19] as a translator, and two classic compression methods: JPEG and BPG as two compressor options for CbT and CaT. In addition, we also train domain-specific neural image compression models on CelebA-HQ and AFHQ separately (NIC in our experiments). We use the same encoder and decoder architecture of StarGAN v2 and a hyperprior entropy model [26] for a fair comparison and optimized with mean square error. The training of NIC is the same with UI2Icodec but without translation mode.

Effect of compression on translation. As shown in Figure 6, we observe changes of both FID and LPIPS values with varying rates (BPP). Notably, CbT always obtains lower FID scores than CaT, which is not surprising since CaT compresses translated images via one lossy codec and the final images have compression artifacts, resulting in worse FID. I2Icodec obtains the lowest FID among all methods at the same rate and also consists of only one encoder and decoder while CaT and CbT are not, which implies that I2Icodec needs less coding time. The diversity measured by LPIPS is shown in Figure 6(b) and Figure 6(d), where CaT can achieve larger scores but with higher distortion, especially for CaT (JPEG) (see the second row in Figure 8). CbT obtains a similar LPIPS as StarGAN v2, which is lower than our I2Icodec. In summary, I2Icodec achieves smaller bandwidth requirements for distributed I2I translation and provides an effective way to guide I2I translation by controlling the amount of information in the bottleneck via the rate constraint. In Table. 1, we also report the quantitative comparison with other I2I translation methods [14, 15, 41, 19] that do not consider compression. It shows that I2Icodec can achieve a range of scores (FID


Figure 7: Translated images with reference on different rate.

Method	Latent-guided synthesis						Reference-guided synthesis					
	CelebA-HQ			AFHQ			CelebA-HQ			AFHQ		
	FID↓	LPIPS↑	BPP	FID↓	LPIPS↑	BPP	FID↓	LPIPS↑	BPP	FID↓	LPIPS↑	BPP
MUNIT [14]	31.4	0.363	-	41.5	0.511	-	107.1	0.176	-	223.9	0.199	-
DRIT [15]	52.1	0.178	-	95.6	0.326	-	53.3	0.311	-	114.8	0.156	-
MSGAN [41]	33.1	0.389	-	61.4	0.517	-	39.6	0.312	-	69.8	0.375	-
StarGANv2 [19]	22.1*	0.115*	-	16.2	0.450	-	23.3*	0.209*	-	19.8	0.432	-
I2Icodec ($\lambda_T = 0.5$)	16.6	0.402	0.043	15.2	0.565	0.042	21.0	0.354	0.043	20.6	0.526	0.042
I2Icodec ($\lambda_T = 0.1$)	20.8	0.224	0.153	14.9	0.533	0.132	20.7	0.247	0.153	20.0	0.494	0.132
I2Icodec ($\lambda_T = 0.05$)	20.2	0.210	0.381	15.4	0.508	0.317	20.0	0.220	0.381	19.9	0.471	0.317
T + A (w GAN)	20.0	0.088	-	25.8	0.288	-	22.2	0.082	-	28.1	0.265	-
T + A (w/o GAN)	20.6	0.098	-	28.7	0.268	-	21.1	0.089	-	32.5	0.238	-
UI2Icodec (T mode)	18.14	0.403	0.065	13.5	0.531	0.131	17.45	0.360	0.065	17.9	0.496	0.131
Real images	14.8	-	-	12.9	-	-	14.8	-	-	12.9	-	-

Table 1

Quantitative comparison. The FIDs of real images are computed between the training and test sets. Note that they may not be optimal values since the number of test images is insufficient, but we report them for reference. * means the results of StarGAN v2 on CelebA-HQ are from the same model architecture on AFHQ, which doesn't include skip connections with the adaptive wing based heatmap [42].

and LPIPS) at different rates for both latent-guided and reference-guided synthesis on two datasets. We want to emphasize that I2Icodec is more efficient and effective than CaT and CbT and also provides a lever for I2I translation.

Visualization of translated images. In Figure 8 we show translated images using different methods. It is obvious that CaT suffers from artifacts (see the second row) even with the better codec BPG (see the third row) and higher rate than other methods. While CaT with NIC as the compressor has less artifacts, the images are blurred (see the fourth row). In addition, different from than CaT with NIC and CbT with NIC, I2Icodec performs both compression and disentanglement jointly with transformation from pixel-level to latent space only once, which is a more efficient way. CbT with JPEG can keep some structure information, but result in unnatural translation (on AFHQ) or serious artifacts (on CelebA-HQ) due to the JPEG compression artifacts themselves. With BPG and NIC, the influence of compression is largely reduced (see the sixth and seventh rows in Figure 8(a)), but note that it still appears again at low

rates (see the sixth row in Figure 8(b)). I2Icodec can generate more natural and diverse images even with extremely low rate. In addition, we also show the synthesized images guided by a reference image on three different rates from high to low in Figure 7. It shows that the translated images have a more similar style to reference image when the rate is lower, and also illustrate that this method can control well how much the translated image obtains the same style of reference image along with the rate.

5.2. Unified I2I translation and autoencoding

In this section, we evaluate the performance of UI2Icodec in I2I translation and autoencoding.

Autoencoding. We compare JPEG, BPG, and NIC as image compression baselines. As shown in Figure 11, UI2Icodec in the autoencoding mode outperforms JPEG largely and BPG marginally on PSNR and MS-SSIM (dB) on similar rates. NIC has the best PSNR results since it was optimized for mean square error. In contrast, UI2Icodec obtains much better LPIPS scores than other methods due

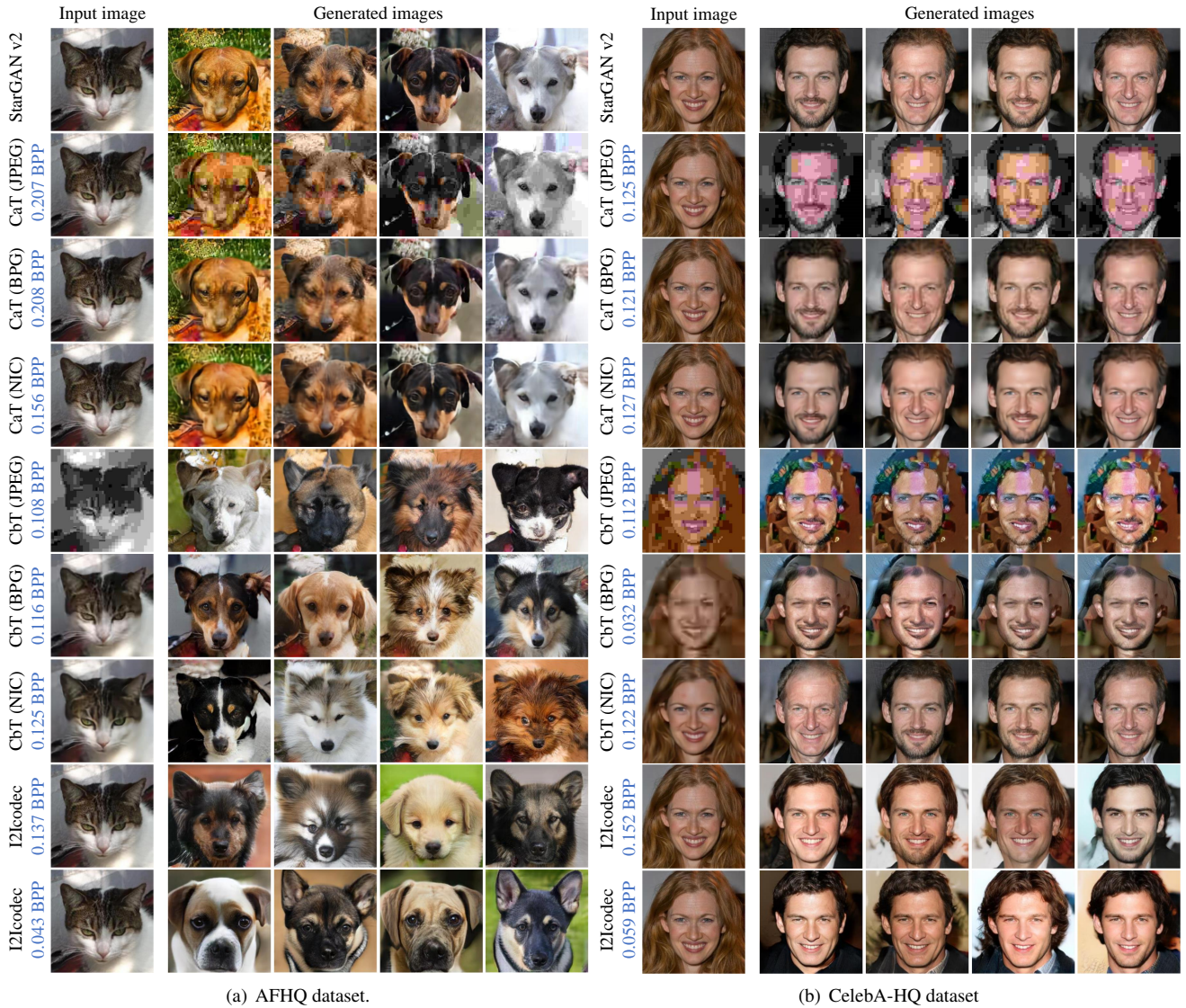


Figure 8: Visualization of latent-guided synthesis with different methods (more samples in Section 5.3).

Method	NIC	StarGAN v2	I2lcodec	UI2lcodec
Number of parameters (millions)	35.30	53.73	54.22	54.23
Training time (hours)	10.3	65.7	71.34	97.6

Table 2

Number of parameters and training time of different methods (CelebA-HQ dataset).

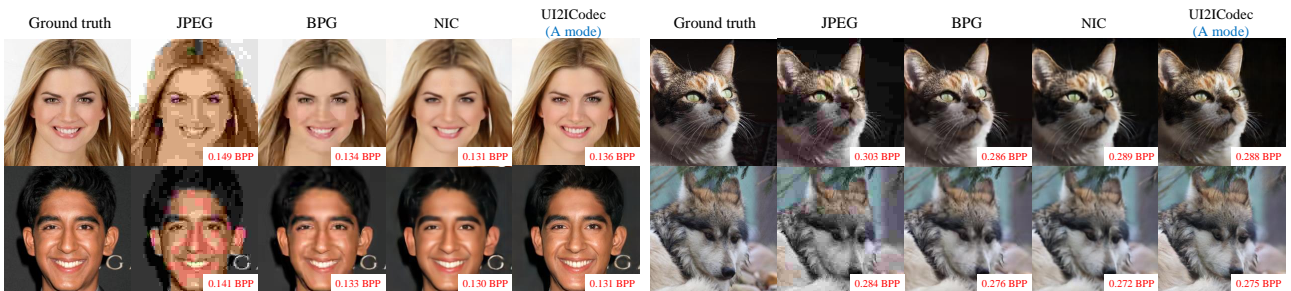


Figure 9: Reconstructions with different compression methods.

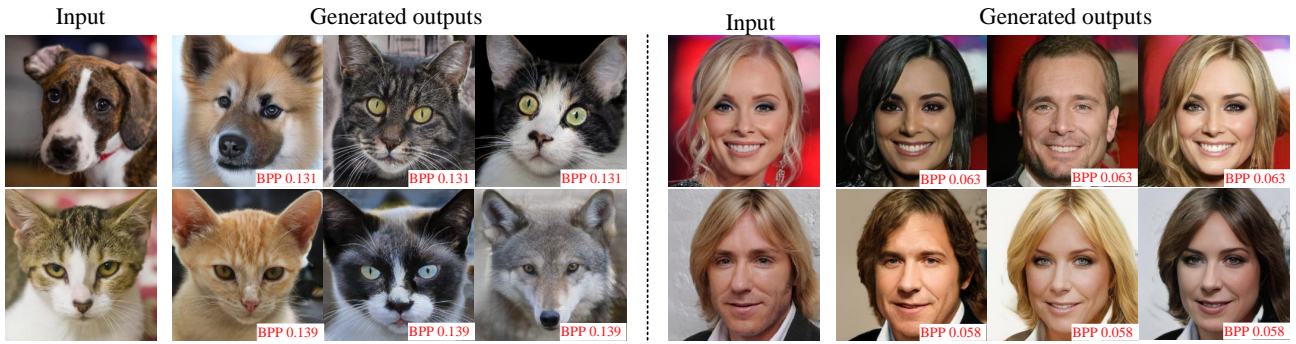


Figure 10: Diverse image synthesis results of UI2Icodec in the translation mode on the CelebA-HQ and AFHQ datasets.

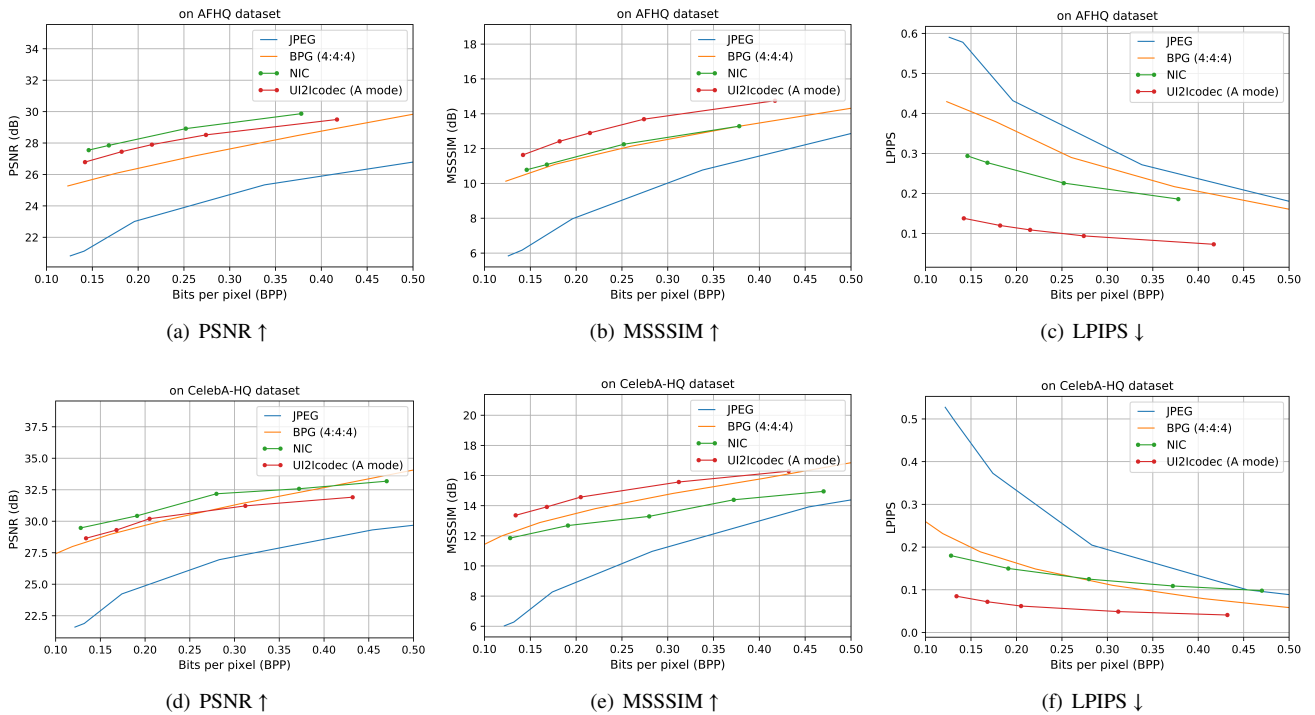


Figure 11: Rate-distortion and -perception curves on AFHQ and CelebA-HQ datasets.

to the adversarial loss. Comparing the examples in Figure 9, we can see that our method can keep more high-frequency information and more natural reconstruction at the similar rate due to the leverage of GAN loss for autoencoding.

I2I translation. The quantitative evaluation of UI2Icodec in the translation mode is shown in Table 1. It shows that it is possible to switch image compression and I2I translation by using our method. In addition, some images generated with UI2Icodec have been already shown in Figure 10. More visualization samples can be viewed in Figure 14 and Figure 15.

5.3. Additional results

Ablation study. We evaluate the effects of the two main modifications of StarGAN v2 architecture: quantization+entropy coding (for compression), and adaptation units (for integrated translation+autoencoding). Comparing

StarGAN v2 and I2Icodec in Table 1, we observe that compression tends to improve FID and LPIPS. In contrast, comparing StarGAN and T+A (StarGAN with adaptation units and autoencoding loss), we observe that the combination of both functionalities has a small penalty in those metrics. However, an important caveat is that FID and LPIPS could be somewhat limited as evaluation metrics in this setting, and further research is required.

User study. In addition to quantitative comparison of I2I translation, we conducted a user study where we asked subjects to select which results they consider more realistic, given the target label and having the same pose as the input image. We apply pairwise comparisons (forced-choice) with 12 users (100 image pairs/user) for I2I translation. Experiments are performed on the AFHQ and CelebA-HQ datasets separately. Figure 12 shows that UI2Icodec (T mode), which runs at the similar BPPs to CaT (BPG) and CbT (BPG),

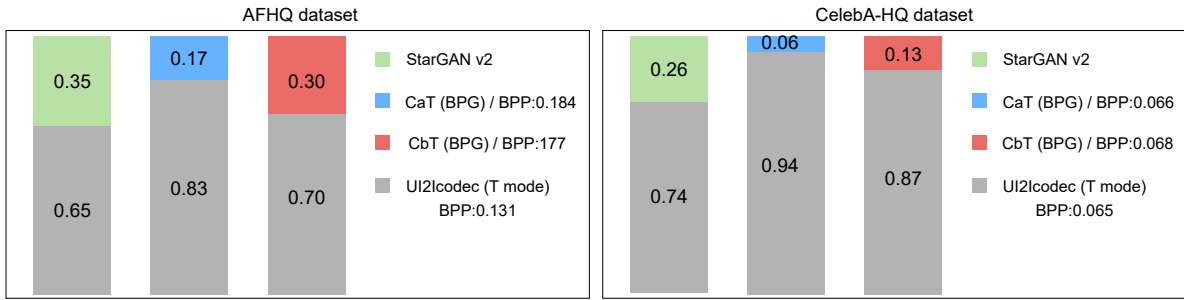


Figure 12: User study for the comparisons between StarGAN v2, CaT (BPG), CbT (BPG) and UI2Icodec (T mode) respectively on AFHQ and CelebA-HQ datasets.

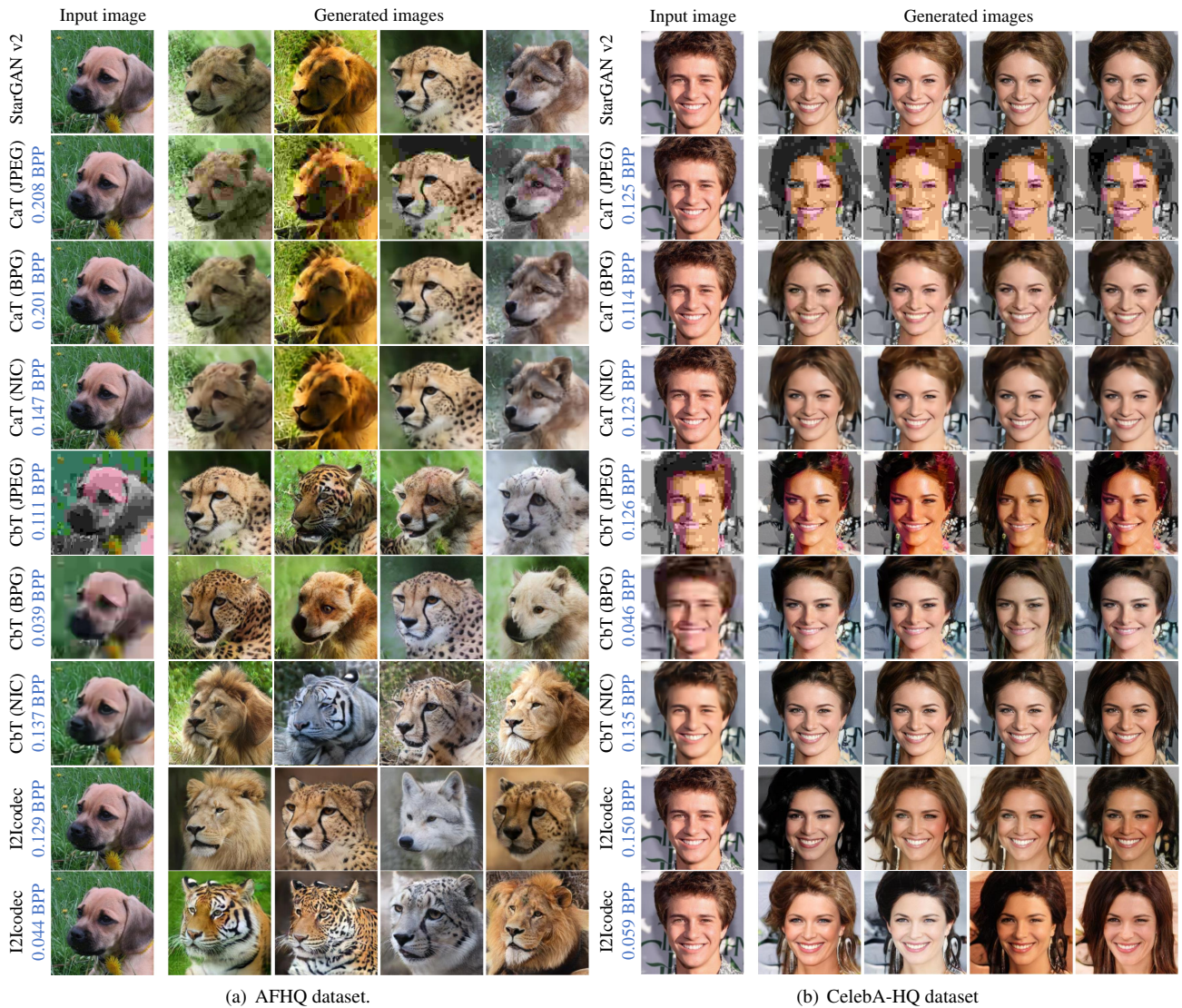


Figure 13: Additional visualization of latent-guided synthesis with different methods.

can obtain better scores than the main baseline methods StarGAN v2, CaT (BPG) and CbT (BPG) respectively.

Model size and training time. Table. 2 shows that the proposed models (both I2Icodec and UI2Icodec) only have 1% more parameters than StarGANv2. Note that CbT and

CaT with NICs require around twice the amount of parameters (since there are two encoders and two decoders). Training requires 8.6%/48% more time (for I2Icodec/UI2Icodec, respectively). Similarly, note that training CbT and CaT with NICs requires training a NIC model and StarGAN v2,

so I2Icodec requires less training time. In addition, we need two independent models when both translation and autoencoding functionalities are needed (i.e. NIC and StarGAN), requiring 89.03M parameters, 65% more than UI2Icodec (54.23M parameters). Finally, the overhead of UI2Icodec with respect to I2Icodec is negligible (only around 0.01%).

More visualization results of I2Icodec and UI2Icodec.

We provide additional latent-guided image synthesis results from two independent I2Icodec models with different rate constraints (see eighth and ninth rows in Figure 13). Same with Figure 8, we also include results of seven baselines: StarGAN v2, Compression after translation (CaT) with JPEG, BPG and NIC, Compression before translation (CbT) with JPEG, BPG and NIC. It shows that CaT and CbT suffer artifacts or result in unnatural or blurred translation, and our I2Icodec can generate natural and diverse images even with an extremely low rate. We also show more results of our UI2Icodec in both translation (T) and autoencoding (A) modes on CelebA-HQ dataset (Figure 14) and AFHQ dataset (Figure 15) separately. It verifies that the UI2Icodec can successfully switch between modes using a single model. Our method can obtain better reconstructions than BPG, both visually and measured in terms of lower LPIPS values.

6. Conclusion

In this paper, we study the novel problem of distributed I2I translation, and its integration with autoencoding in a joint model, resulting in the proposed I2Icodec and UI2Icodec frameworks. Distributed I2I translation required augmenting an I2I translation framework with quantization and entropy coding. Interestingly, constraining the rate can control the amount of source information in I2I translation. The experiments show that our joint model can keep competitive autoencoding and translation performance.

7. Acknowledgments

We acknowledge the support from Huawei Kirin Solution and the Spanish Government funding for projects RTI2018-102285-A-I00 and RYC2019-027020-I.

References

- [1] F. Bellard, BPG Image Format, <http://bellard.org/bpg/>, [Online; accessed 8-January-2020] (2014).
- [2] P. Isola, J.-Y. Zhu, T. Zhou, A. A. Efros, Image-to-image translation with conditional adversarial networks, 2017.
- [3] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, E. Shechtman, Toward multimodal image-to-image translation, 2017, pp. 465–476.
- [4] A. Gonzalez-Garcia, J. van de Weijer, Y. Bengio, Image-to-image translation for cross-domain disentanglement, 2018, pp. 1294–1305.
- [5] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, B. Catanzaro, Video-to-video synthesis, 2018.
- [6] T.-C. Wang, M.-Y. Liu, A. Tao, G. Liu, J. Kautz, B. Catanzaro, Few-shot video-to-video synthesis, 2019.
- [7] M.-Y. Liu, T. Breuel, J. Kautz, Unsupervised image-to-image translation networks, 2017, pp. 700–708.
- [8] T. Kim, M. Cha, H. Kim, J. Lee, J. Kim, Learning to discover cross-domain relations with generative adversarial networks, 2017.
- [9] Z. Yi, H. Zhang, P. T. Gong, et al., Dualgan: Unsupervised dual learning for image-to-image translation, 2017.
- [10] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, 2017, pp. 2223–2232.
- [11] Z. Gan, L. Chen, W. Wang, Y. Pu, Y. Zhang, H. Liu, C. Li, L. Carin, Triangle generative adversarial networks, 2017, pp. 5247–5256.
- [12] Y. A. Mejjati, C. Richardt, J. Tompkin, D. Cosker, K. I. Kim, Unsupervised attention-guided image-to-image translation, 2018, pp. 3693–3703.
- [13] T. Park, A. A. Efros, R. Zhang, J.-Y. Zhu, Contrastive learning for conditional image synthesis, 2020.
- [14] X. Huang, M.-Y. Liu, S. Belongie, J. Kautz, Multimodal unsupervised image-to-image translation, 2018, pp. 172–189.
- [15] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. K. Singh, M.-H. Yang, Diverse image-to-image translation via disentangled representations, 2018.
- [16] A. Almahairi, S. Rajeswar, A. Sordoni, P. Bachman, A. Courville, Augmented cyclegan: Learning many-to-many mappings from unpaired data, 2018.
- [17] A. Royer, K. Bousmalis, S. Gouws, F. Bertsch, I. Mosseri, F. Cole, K. Murphy, Xgan: Unsupervised image-to-image translation for many-to-many mappings, 2018.
- [18] X. Yu, Y. Chen, S. Liu, T. Li, G. Li, Multi-mapping image-to-image translation via learning disentanglement, 2019, pp. 2990–2999.
- [19] Y. Choi, Y. Uh, J. Yoo, J.-W. Ha, Stargan v2: Diverse image synthesis for multiple domains, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 8188–8197.
- [20] H.-Y. Lee, H.-Y. Tseng, Q. Mao, J.-B. Huang, Y.-D. Lu, M. Singh, M.-H. Yang, Drit++: Diverse image-to-image translation via disentangled representations (2020) 1–16.
- [21] G. Toderici, S. M. O’Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, R. Sukthankar, Variable rate image compression with recurrent neural networks, arXiv preprint arXiv:1511.06085.
- [22] J. Ballé, V. Laparra, E. P. Simoncelli, End-to-end optimized image compression, arXiv preprint arXiv:1611.01704.
- [23] L. Theis, W. Shi, A. Cunningham, F. Huszár, Lossy image compression with compressive autoencoders, arXiv preprint arXiv:1703.00395.
- [24] E. Agustsson, F. Mentzer, M. Tschannen, L. Cavigelli, R. Timofte, L. Benini, L. V. Gool, Soft-to-hard vector quantization for end-to-end learning compressible representations, in: Advances in Neural Information Processing Systems, 2017, pp. 1141–1151.
- [25] G. Toderici, D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, M. Covell, Full resolution image compression with recurrent neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 5306–5314.
- [26] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, N. Johnston, Variational image compression with a scale hyperprior, in: International Conference on Learning Representations, 2018.
- [27] D. Minnen, J. Ballé, G. D. Toderici, Joint autoregressive and hierarchical priors for learned image compression, in: Advances in Neural Information Processing Systems, 2018, pp. 10771–10780.
- [28] J. Lee, S. Cho, S.-K. Beack, Context-adaptive entropy model for end-to-end optimized image compression, in: International Conference on Learning Representations, 2018.
- [29] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, L. Van Gool, Conditional probability models for deep image compression, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 4394–4402.
- [30] M. Li, K. Ma, J. You, D. Zhang, W. Zuo, Efficient and effective context-based convolutional entropy modeling for image compression, IEEE Transactions on Image Processing 29 (2020) 5900–5911.
- [31] D. Minnen, S. Singh, Channel-wise autoregressive entropy models for learned image compression, in: 2020 IEEE International Conference on Image Processing (ICIP), IEEE, 2020, pp. 3339–3343.

- [32] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, 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.
- [33] Y. Blau, T. Michaeli, Rethinking lossy compression: The rate-distortion-perception tradeoff, in: International Conference on Machine Learning, PMLR, 2019, pp. 675–685.
- [34] Y. Choi, M. El-Khamy, J. Lee, Variable rate deep image compression with a conditional autoencoder, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 3146–3154.
- [35] F. Yang, L. Herranz, J. van de Weijer, J. Iglesias-Guitián, A. López, M. Mozerov, Variable rate deep image compression with modulated autoencoder, *IEEE Signal Processing Letters* 27 (2020) 331–335.
- [36] Y. Patel, S. Appalaraju, R. Manmatha, Deep perceptual compression, arXiv preprint arXiv:1907.08310.
- [37] Y. Patel, S. Appalaraju, R. Manmatha, Hierarchical auto-regressive model for image compression incorporating object saliency and a deep perceptual loss, arXiv preprint arXiv:2002.04988.
- [38] Y. Patel, S. Appalaraju, R. Manmatha, Saliency driven perceptual image compression, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2021, pp. 227–236.
- [39] F. Mentzer, G. D. Toderici, M. Tschannen, E. Agustsson, High-fidelity generative image compression, *Advances in Neural Information Processing Systems* 33 (2020) 11913–11924.
- [40] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, 2014, pp. 2672–2680.
- [41] Q. Mao, H.-Y. Lee, H.-Y. Tseng, S. Ma, M.-H. Yang, Mode seeking generative adversarial networks for diverse image synthesis, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 1429–1437.
- [42] X. Wang, L. Bo, L. Fuxin, Adaptive wimg loss for robust face alignment via heatmap regression, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 6971–6981.
- [43] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.
- [44] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, S. Hochreiter, Gans trained by a two time-scale update rule converge to a local nash equilibrium, *Advances in neural information processing systems* 30.
- [45] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, 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.

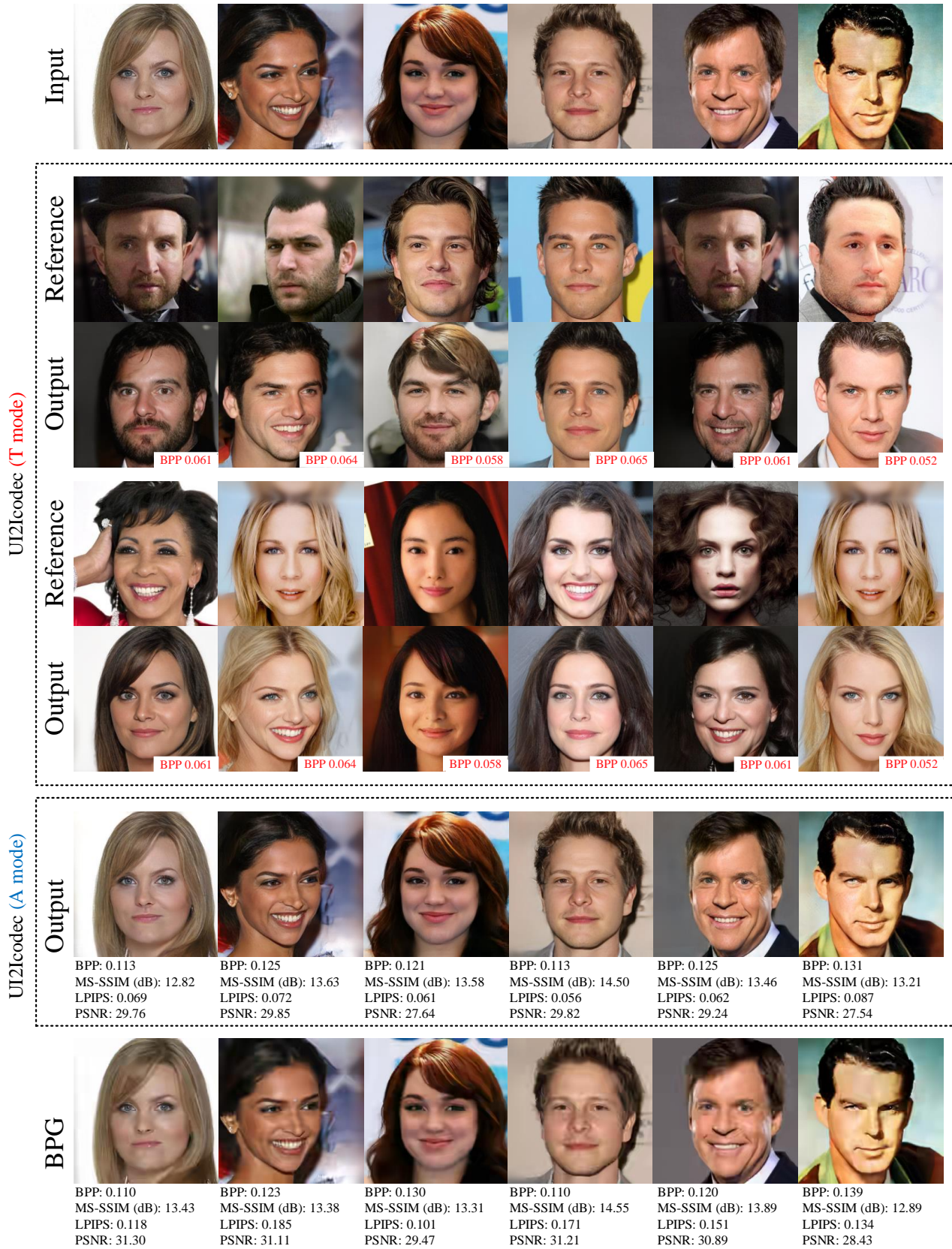


Figure 14: Diverse image synthesis results of UI2Icodec (reference-guided) in the translation and autoencoding modes on CelebA-HQ dataset.

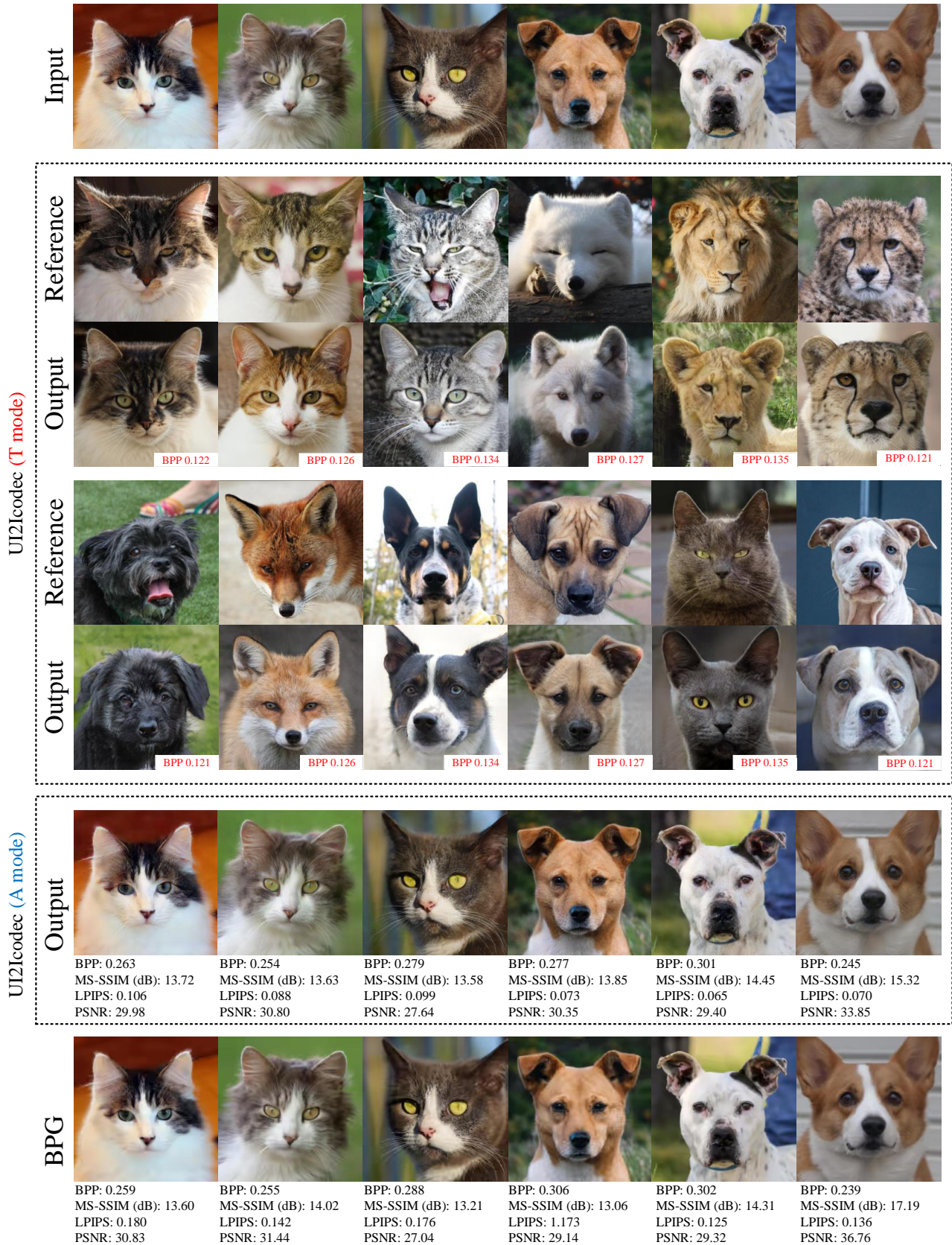


Figure 15: Diverse image synthesis results of UI2Icodec (reference-guided) in the translation and autoencoding modes on AFHQ dataset.