

# DRIT++: Diverse Image-to-Image Translation via Disentangled Representations

Hsin-Ying Lee\* · Hung-Yu Tseng\* · Qi Mao\* · Jia-Bin Huang ·  
Yu-Ding Lu · Maneesh Singh · Ming-Hsuan Yang

Received: date / Accepted: date

**Abstract** Image-to-image translation aims to learn the mapping between two visual domains. There are two main challenges for this task: 1) lack of aligned training pairs and 2) multiple possible outputs from a single input image. In this work, we present an approach based on disentangled representation for generating diverse outputs without paired training images. To synthesize diverse outputs, we propose to embed images onto two spaces: a domain-invariant content space capturing shared information across domains and a domain-specific attribute space. Our model takes the encoded content features extracted from a given input and attribute vectors sampled from the attribute space to synthesize diverse outputs at test time. To handle unpaired training data, we introduce a cross-cycle consistency loss based on disentangled representations. Qualitative results show that our model can generate diverse and realistic images on a wide range of tasks without paired training data. For quantitative evaluations, we

\* Equal contribution

Hsin-Ying Lee, Hung-Yu Tseng, Yu-Ding Lu, and Ming-Hsuan Yang  
Electrical Engineering and Computer Science, University of California at Merced, Merced, CA 95343  
E-mail: {hlee246, htseng6, ylu52, mhyang}@ucmerced.edu

Qi Mao  
Electrical Engineering and Computer Science, Peking University, Beijing, China  
E-mail: qimao@pku.edu.cn

Jia-Bin Huang  
Electrical and Computer Engineering, Virginia Tech, Blacksburg, VA 24060  
E-mail: jhuang@vt.edu

Maneesh Singh  
Verisk Analytics, Jersey City, NJ 07310  
E-mail: maneesh.singh@verisk.com

measure realism with user study and Fréchet inception distance, and measure diversity with the perceptual distance metric, Jensen-Shannon divergence, and number of statistically-different bins.

## 1 Introduction

Image-to-Image (I2I) translation aims to learn the mapping between different visual domains. Numerous vision and graphics problems can be formulated as I2I translation problems, such as colorization [22, 47] (grayscale → color), super-resolution [23, 21, 25] (low-resolution → high-resolution), and photorealistic image synthesis [5, 44, 34] (label → image). In addition, I2I translation can be applied to synthesize images for domain adaptation [3, 40, 15, 33, 7].

Learning the mapping between two visual domains is challenging for two main reasons. First, aligned training image pairs are either difficult to collect (e.g., day scene ↔ night scene) or do not exist (e.g., artwork ↔ real photo). Second, many such mappings are inherently multimodal — a single input may correspond to multiple possible outputs. To handle multimodal translation, one possible approach is to inject a random noise vector to the generator for modeling the multimodal data distribution in the target domain. However, mode collapse may still occur easily since the generator often ignores the additional noise vectors.

Several recent efforts have been made to address these issues. The Pix2pix [17] method applies conditional generative adversarial network to I2I translation problems. Nevertheless, the training process requires paired data. A number of recent approaches [49, 27, 45, 42, 9] relax the dependency on paired training data for learning I2I translation. These methods, however, generate a single output conditioned on the given input

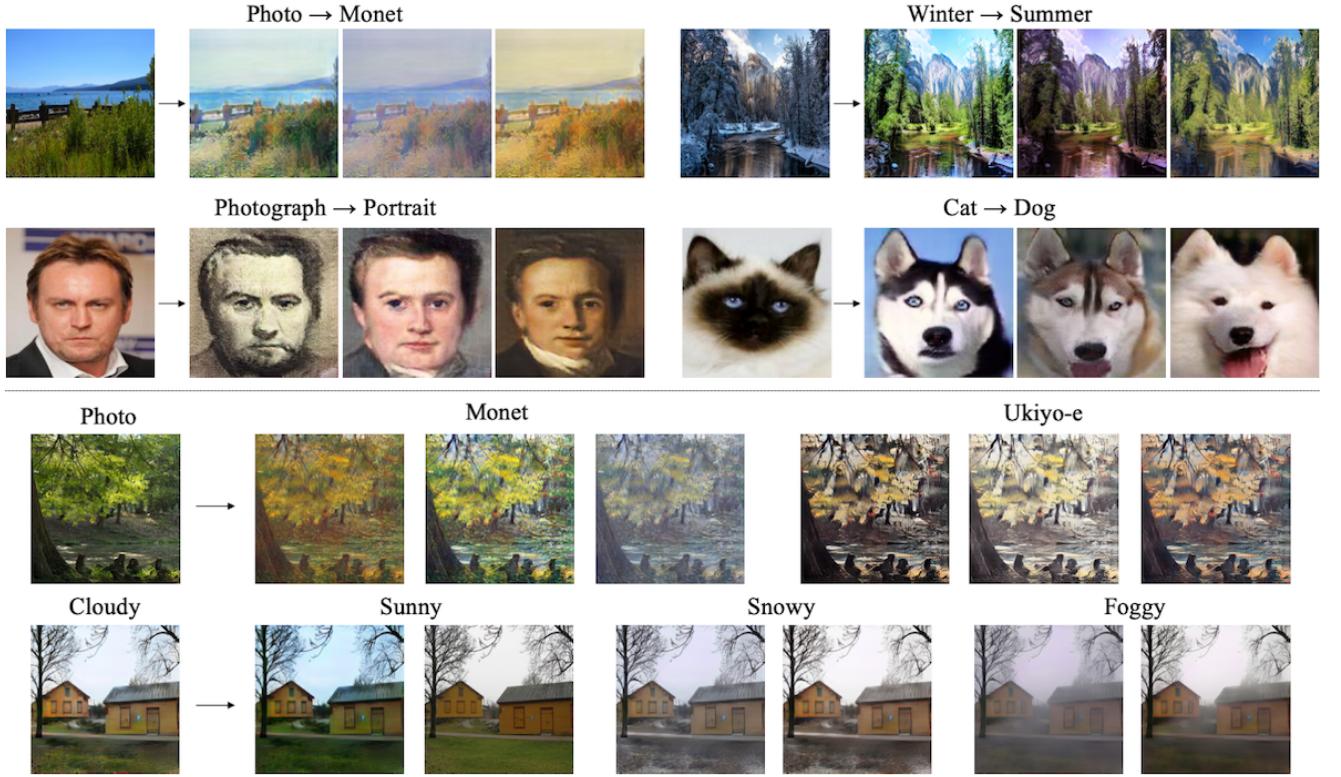


Fig. 1: **Unpaired diverse image-to-image translation.** (Top) Our model learns to perform diverse translation between two collections of images without aligned training pairs. (Bottom) Multi-domain image-to-image translation.

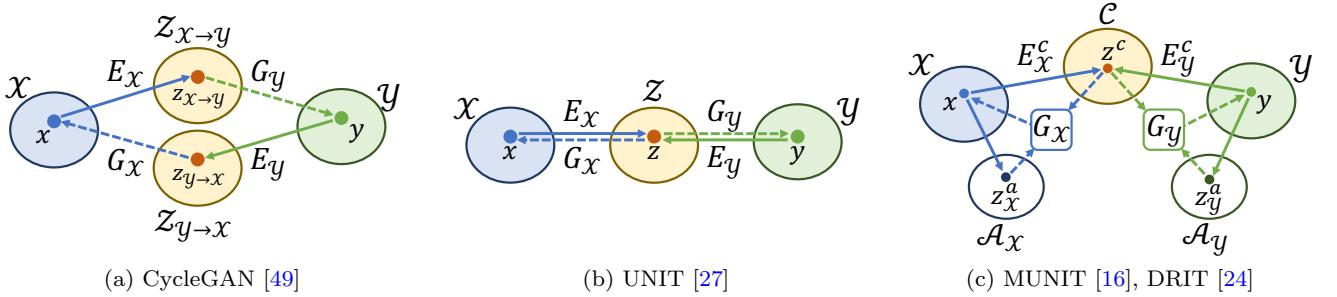


Fig. 2: **Comparisons of unsupervised I2I translation methods.** Denote  $x$  and  $y$  as images in domain  $\mathcal{X}$  and  $\mathcal{Y}$ : (a) CycleGAN [49] maps  $x$  and  $y$  onto *separated* latent spaces. (b) UNIT [27] assumes  $x$  and  $y$  can be mapped onto a *shared* latent space. (c) Our approach disentangles the latent spaces of  $x$  and  $y$  into a shared content space  $\mathcal{C}$  and an attribute space  $\mathcal{A}$  of each domain.

image. As shown in [17, 50], the strategy of incorporating noise vectors as additional inputs to the generator does not increase variations of generated outputs due to the mode collapse issue. The generators in these methods are likely to overlook the added noise vectors. Most recently, the BicycleGAN [50] algorithm tackles the problem of generating diverse outputs in I2I translation by encouraging the one-to-one relationship between the output and the latent vector. Nevertheless,

the training process of BicycleGAN requires paired images.

In this paper, we propose a disentangled representation framework for learning to generate *diverse* outputs with *unpaired* training data. We propose to embed images onto two spaces: 1) a domain-invariant content space and 2) a domain-specific attribute space as shown in Figure 2. Our generator learns to perform I2I translation conditioned on content features and a latent attribute vector. The domain-specific attribute

space aims to model variations within a domain given the same content, while the domain-invariant content space captures information across domains. We disentangle the representations by applying a content adversarial loss to encourage the content features *not* to carry domain-specific cues, and a latent regression loss to encourage the invertible mapping between the latent attribute vectors and the corresponding outputs. To handle unpaired datasets, we propose a *cross-cycle consistency loss* using the proposed disentangled representations. Given a pair of unaligned images, we first perform a cross-domain mapping to obtain intermediate results by swapping the attribute vectors from both images. We can then reconstruct the original input image pair by applying the cross-domain mapping one more time and use the proposed cross-cycle consistency loss to enforce the consistency between the original and the reconstructed images. Furthermore, we apply the mode seeking regularization [31] to further improve the diversity of generated images. At test time, we can use either 1) randomly sampled vectors from the attribute space to generate diverse outputs or 2) the transferred attribute vectors extracted from existing images for example-guided translation. Figure 1 shows examples of diverse outputs produced by our model.

We evaluate the proposed model with extensive qualitative and quantitative experiments. For various I2I tasks, we show diverse translation results with randomly sampled attribute vectors and example-guided translation with transferred attribute vectors from existing images. In addition to the common two-domain image-to-image translation, we extend our proposed framework to the more general multi-domain image-to-image translation and demonstrate diverse translation among domains. We measure realism of our results with a user study and the Fréchet inception distance (FID) [14], and evaluate diversity using perceptual distance metrics [48]. However, the diversity metric alone does not effectively measure similarity between the distribution of generated images and the distribution of real data. Therefore, we use the Jensen-Shannon Divergence (JSD) distance which measures the similarity between distributions, and the Number of Statistically-Different Bins (NDB) [38] metric which determines the relative proportions of samples within clusters predetermined by real data.

We make the following contributions in this work:

- 1) We introduce a disentangled representation framework for image-to-image translation. We apply a content discriminator to facilitate the factorization of domain-invariant content space and domain-specific attribute space, and a cross-cycle consistency loss that allows us to train the model with unpaired data.

- 2) Extensive qualitative and quantitative experiments show that our model performs favorably against existing I2I models. Images generated by our model are both diverse and realistic.

- 3) The proposed disentangled representation and cross-cycle consistency can be applied to multi-domain image-to-image translation for generating diverse images.

## 2 Related Work

**Generative adversarial networks.** The recent years have witnessed rapid advances of generative adversarial networks (GANs) [13, 36, 2] for image generation. The core idea of GANs lies in the adversarial loss that enforces the distribution of generated images to match that of the target domain. The generators in GANs can map from noise vectors to realistic images. Several recent efforts exploit *conditional* GAN in various contexts including conditioned on text [37], low-resolution images [23], video frames [43], and image [17]. Our work focuses on using GAN conditioned on an input image. In contrast to several existing conditional GAN frameworks that require paired training data, our model generates diverse outputs without paired data. As such, our method has wider applicability to problems where paired training datasets are scarce or not available.

**Image-to-image translation.** I2I translation aims to learn the mapping from a source image domain to a target image domain. The Pix2pix [17] method applies a conditional GAN to model the mapping function. Although high-quality results have been shown, the model training requires paired training data. To train with unpaired data, the CycleGAN [49], DiscoGAN [18], and UNIT [27] schemes leverage cycle consistency to regularize the training. However, these methods perform generation conditioned solely on an input image and thus produce one single output. Simply injecting a noise vector to a generator is usually not an effective solution to achieve multimodal generation due to the lack of regularization between the noise vectors and the target domain. On the other hand, the BicycleGAN [50] algorithm enforces the bijection mapping between the latent and target space to tackle the mode collapse problem. Nevertheless, the method is only applicable to problems with paired training data. Unlike existing work, our method enables I2I translation with diverse outputs in the absence of paired training data.

We note several concurrent methods [1, 16, 4, 29] (all independently developed) also adopt disentangled representations similar to our work for learning diverse I2I translation from unpaired training data. Furthermore,

several approaches [10, 26] extend the conventional two-domain I2I to general multi-domain settings. However, these methods can only achieve one-to-one mapping among domains.

**Disentangled representations.** The task of learning disentangled representation aims at modeling the factors of data variations. Previous work makes use of labeled data to factorize representations into class-related and class-independent components [8, 20, 30, 32]. Recently, numerous unsupervised methods have been developed [6, 12] to learn disentangled representations. The InfoGAN [6] algorithm achieves disentanglement by maximizing the mutual information between latent variables and data variation. Similar to DrNet [12] that separates time-independent and time-varying components with an adversarial loss, we apply a content adversarial loss to disentangle an image into domain-invariant and domain-specific representations to facilitate learning diverse cross-domain mappings.

### 3 Disentangled Representation for I2I Translation

Our goal is to learn a multimodal mapping between two visual domains  $\mathcal{X} \subset \mathbb{R}^{H \times W \times 3}$  and  $\mathcal{Y} \subset \mathbb{R}^{H \times W \times 3}$  without paired training data. As illustrated in Figure 3, our framework consists of content encoders  $\{E_{\mathcal{X}}^c, E_{\mathcal{Y}}^c\}$ , attribute encoders  $\{E_{\mathcal{X}}^a, E_{\mathcal{Y}}^a\}$ , generators  $\{G_{\mathcal{X}}, G_{\mathcal{Y}}\}$ , and domain discriminators  $\{D_{\mathcal{X}}, D_{\mathcal{Y}}\}$  for both domains, and a content discriminators  $D_{\text{adv}}^c$ . Taking domain  $\mathcal{X}$  as an example, the content encoder  $E_{\mathcal{X}}^c$  maps images onto a shared, domain-invariant content space ( $E_{\mathcal{X}}^c : \mathcal{X} \rightarrow \mathcal{C}$ ) and the attribute encoder  $E_{\mathcal{X}}^a$  maps images onto a domain-specific attribute space ( $E_{\mathcal{X}}^a : \mathcal{X} \rightarrow \mathcal{A}_{\mathcal{X}}$ ). The generator  $G_{\mathcal{X}}$  synthesizes images conditioned on both content and attribute vectors ( $G_{\mathcal{X}} : \{\mathcal{C}, \mathcal{A}_{\mathcal{X}}\} \rightarrow \mathcal{X}$ ). The discriminator  $D_{\mathcal{X}}$  aims to discriminate between real images and translated images in the domain  $\mathcal{X}$ . In addition, the content discriminator  $D^c$  is trained to distinguish the extracted content representations between two domains. To synthesize multimodal outputs at test time, we regularize the attribute vectors so that they can be drawn from a prior Gaussian distribution  $N(0, 1)$ .

#### 3.1 Disentangle Content and Attribute Representations

Our approach embeds input images onto a shared content space  $\mathcal{C}$ , and domain-specific attribute spaces,  $\mathcal{A}_{\mathcal{X}}$  and  $\mathcal{A}_{\mathcal{Y}}$ . Intuitively, the content encoders should encode the common information that is *shared* between

domains onto  $\mathcal{C}$ , while the attribute encoders should map the remaining domain-specific information onto  $\mathcal{A}_{\mathcal{X}}$  and  $\mathcal{A}_{\mathcal{Y}}$ .

$$\begin{aligned} \{z_x^c, z_x^a\} &= \{E_{\mathcal{X}}^c(x), E_{\mathcal{X}}^a(x)\} & z_x^c \in \mathcal{C}, z_x^a \in \mathcal{A}_{\mathcal{X}} \\ \{z_y^c, z_y^a\} &= \{E_{\mathcal{Y}}^c(y), E_{\mathcal{Y}}^a(y)\} & z_y^c \in \mathcal{C}, z_y^a \in \mathcal{A}_{\mathcal{Y}} \end{aligned} \quad (1)$$

To achieve representation disentanglement, we apply two strategies: weight-sharing and a content discriminator. First, similar to [27], based on the assumption that two domains share a common latent space, we share the weight between the last layer of  $E_{\mathcal{X}}^c$  and  $E_{\mathcal{Y}}^c$  and the first layer of  $G_{\mathcal{X}}$  and  $G_{\mathcal{Y}}$ . Through weight sharing, we force the content representation to be mapped onto the same space. However, sharing the same high-level mapping functions cannot guarantee the same content representations encode the same information for both domains. Therefore, we propose a content discriminator  $D^c$  which aims to distinguish the domain membership of the encoded content features  $z_x^c$  and  $z_y^c$ . On the other hand, content encoders learn to produce encoded content representations whose domain membership cannot be distinguished by the content discriminator  $D^c$ . We express this content adversarial loss as:

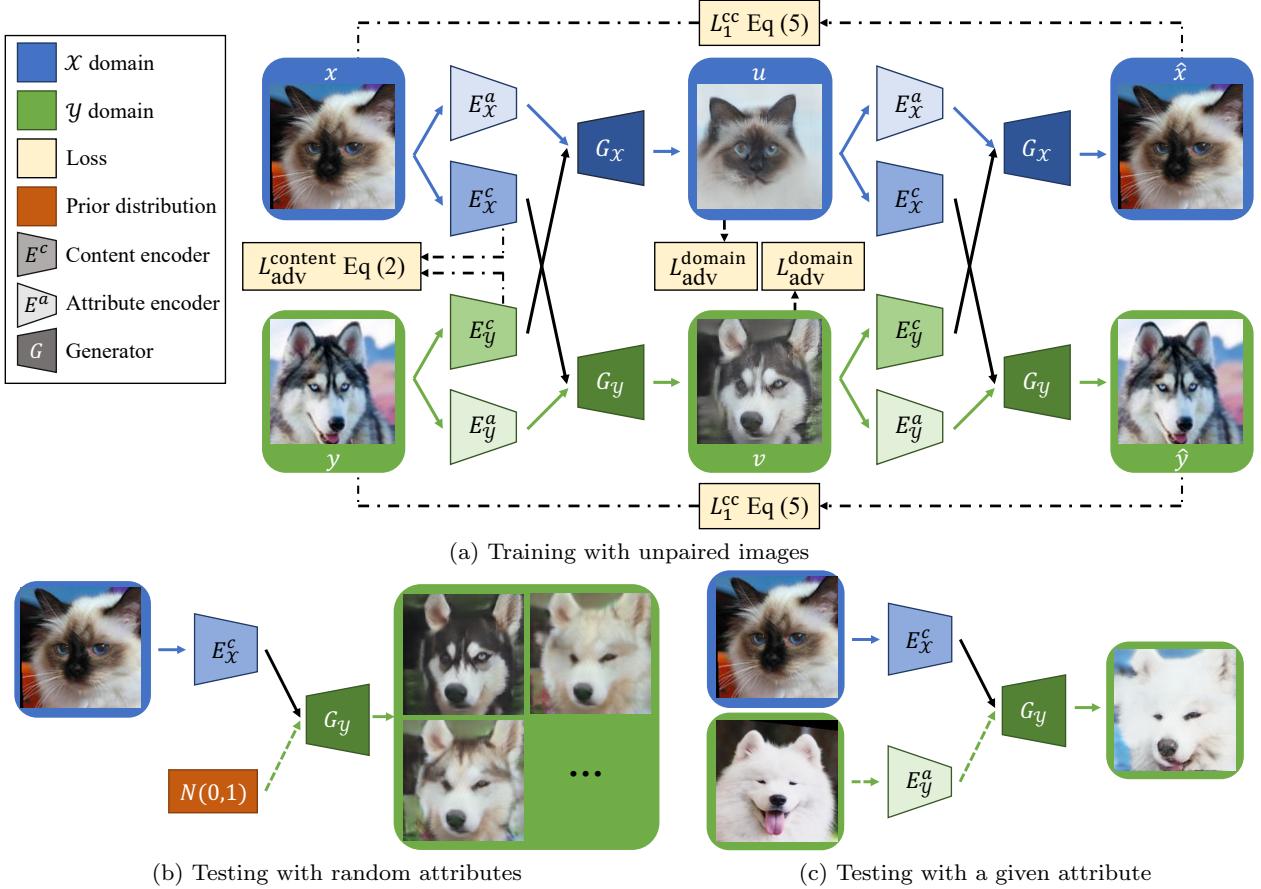
$$\begin{aligned} L_{\text{adv}}^{\text{content}}(E_{\mathcal{X}}^c, E_{\mathcal{Y}}^c, D^c) = & \mathbb{E}_x [\frac{1}{2} \log D^c(E_{\mathcal{X}}^c(x)) + \frac{1}{2} \log (1 - D^c(E_{\mathcal{X}}^c(x)))] \\ & + \mathbb{E}_y [\frac{1}{2} \log D^c(E_{\mathcal{Y}}^c(y)) + \frac{1}{2} \log (1 - D^c(E_{\mathcal{Y}}^c(y)))] \end{aligned} \quad (2)$$

#### 3.2 Cross-cycle Consistency Loss

With the disentangled representation where the content space is shared among domains and the attribute space encodes intra-domain variations, we can perform I2I translation by combining a content representation from an arbitrary image and an attribute representation from an image of the target domain. We leverage this property and propose a *cross-cycle consistency*. In contrast to cycle consistency constraint in [49] (i.e.,  $\mathcal{X} \rightarrow \mathcal{Y} \rightarrow \mathcal{X}$ ) which assumes one-to-one mapping between the two domains, the proposed cross-cycle constraint exploit the disentangled content and attribute representations for cyclic reconstruction.

Our cross-cycle constraint consists of two stages of I2I translation.

**Forward translation.** Given a non-corresponding pair of images  $x$  and  $y$ , we encode them into  $\{z_x^c, z_x^a\}$  and  $\{z_y^c, z_y^a\}$ . We then perform the first translation by swapping the attribute representation (i.e.,  $z_x^a$  and  $z_y^a$ ) to generate  $\{u, v\}$ , where  $u \in \mathcal{X}, v \in \mathcal{Y}$ .



**Fig. 3: Method overview.** (a) With the proposed content adversarial loss  $L_{\text{adv}}^{\text{content}}$  (Section 3.1) and the cross-cycle consistency loss  $L_1^{\text{cc}}$  (Section 3.2), we are able to learn the multimodal mapping between the domain  $\mathcal{X}$  and  $\mathcal{Y}$  with unpaired data. Thanks to the proposed disentangled representation, we can generate output images conditioned on either (b) random attributes or (c) a given attribute at test time.

$$u = G_{\mathcal{X}}(z_y^c, z_x^a) \quad v = G_{\mathcal{Y}}(z_x^c, z_y^a) \quad (3)$$

### 3.3 Other Loss Functions

**Backward translation.** After encoding  $u$  and  $v$  into  $\{z_u^c, z_u^a\}$  and  $\{z_v^c, z_v^a\}$ , we perform the second translation by once again swapping the attribute representation (i.e.,  $z_u^a$  and  $z_v^a$ ).

$$\hat{x} = G_{\mathcal{X}}(z_v^c, z_u^a) \quad \hat{y} = G_{\mathcal{Y}}(z_u^c, z_v^a) \quad (4)$$

Here, after two I2I translation stages, the translation should reconstruct the original images  $x$  and  $y$  (as illustrated in Figure 3). To enforce this constraint, we formulate the *cross-cycle consistency loss* as:

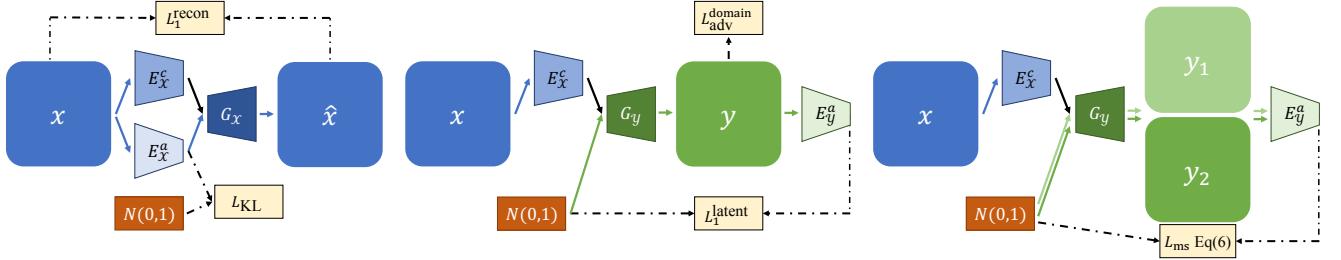
$$\begin{aligned} L_1^{\text{cc}}(G_{\mathcal{X}}, G_{\mathcal{Y}}, E_{\mathcal{X}}^c, E_{\mathcal{Y}}^c, E_{\mathcal{X}}^a, E_{\mathcal{Y}}^a) = \\ \mathbb{E}_{x,y} [\|G_{\mathcal{X}}(E_{\mathcal{Y}}^c(v), E_{\mathcal{X}}^a(u)) - x\|_1 \\ + \|G_{\mathcal{Y}}(E_{\mathcal{X}}^c(u), E_{\mathcal{Y}}^a(v)) - y\|_1], \end{aligned} \quad (5)$$

where  $u = G_{\mathcal{X}}(E_{\mathcal{Y}}^c(y), E_{\mathcal{X}}^a(x))$  and  $v = G_{\mathcal{Y}}(E_{\mathcal{X}}^c(x), E_{\mathcal{Y}}^a(y))$ .  $\hat{y} = G_{\mathcal{Y}}(E_{\mathcal{Y}}^c(y), E_{\mathcal{X}}^a(x))$ .

Other than the proposed content adversarial loss and cross-cycle consistency loss, we also use several other loss functions to facilitate network training. We illustrate these additional losses in Figure 4. Starting from the top-right, in the counter-clockwise order:

**Domain adversarial loss.** We impose adversarial loss  $L_{\text{adv}}^{\text{domain}}$  where  $D_{\mathcal{X}}$  and  $D_{\mathcal{Y}}$  attempt to discriminate between real images and generated images in each domain, while  $G_{\mathcal{X}}$  and  $G_{\mathcal{Y}}$  attempt to generate realistic images.

**Self-reconstruction loss.** In addition to the cross-cycle reconstruction, we apply a self-reconstruction loss  $L_1^{\text{rec}}$  to facilitate the training. With encoded content and attribute features  $\{z_x^c, z_x^a\}$  and  $\{z_y^c, z_y^a\}$ , the decoders  $G_{\mathcal{X}}$  and  $G_{\mathcal{Y}}$  should decode them back to original input  $x$  and  $y$ . That is,  $\hat{x} = G_{\mathcal{X}}(E_{\mathcal{X}}^c(x), E_{\mathcal{X}}^a(x))$  and  $\hat{y} = G_{\mathcal{Y}}(E_{\mathcal{Y}}^c(y), E_{\mathcal{Y}}^a(y))$ .



**Fig. 4: Additional loss functions.** In addition to the cross-cycle reconstruction loss  $L_1^{cc}$  and the content adversarial loss  $L_{adv}^{content}$  described in Figure 3, we apply several additional loss functions in our training process. The self-reconstruction loss  $L_1^{recon}$  facilitates training with self-reconstruction; the KL loss  $L_{KL}$  aims to align the attribute representation with a prior Gaussian distribution; the adversarial loss  $L_1^{domain}$  encourages  $G$  to generate realistic images in each domain; and the latent regression loss  $L_1^{latent}$  enforces the reconstruction on the latent attribute vector. Finally, the mode seeking regularization  $L_{ms}$  further improves the diversity. More details can be found in Section 3.3 and Section 3.4.

**KL loss.** In order to perform stochastic sampling at test time, we encourage the attribute representation to be as close to a prior Gaussian distribution. We thus apply the loss  $L_{KL} = \mathbb{E}[D_{KL}((z_a) \| N(0, 1))]$ , where  $D_{KL}(p\|q) = -\int p(z) \log \frac{p(z)}{q(z)} dz$ .

**Latent regression loss.** To encourage invertible mapping between the image and the latent space, we apply a latent regression loss  $L_1^{latent}$  similar to [50]. We draw a latent vector  $z$  from the prior Gaussian distribution as the attribute representation and attempt to reconstruct it with  $\hat{z} = E_X^a(G_X(E_X^c(x), z))$  and  $\hat{z} = E_Y^a(G_Y(E_Y^c(y), z))$ .

The full objective function of our network is:

$$\begin{aligned} L_{D,D^c} &= \lambda_{adv}^{content} L_{adv}^c + \lambda_{adv}^{domain} L_{adv}^d \\ L_{G,E^c,E^a} &= -L_{D,D^c} + \lambda_1^{cc} L_1^{cc} + \lambda_1^{recon} L_1^{recon} \\ &\quad + \lambda_1^{latent} L_1^{latent} + \lambda_{KL} L_{KL} \end{aligned}$$

where the hyper-parameters  $\lambda$ s control the importance of each term.

### 3.4 Mode Seeking Regularization

We incorporate the mode seeking regularization [31] alleviate the mode-dropping problem in conditional generation tasks. Given a conditional image  $\mathbf{I}$ , latent vectors  $\mathbf{z}_1$  and  $\mathbf{z}_2$ , and a conditional generator  $G$ , the mode seeking regularization term aims to maximize the ratio of the distance between  $G(\mathbf{I}, \mathbf{z}_1)$  and  $G(\mathbf{I}, \mathbf{z}_2)$  with respect to the distance between  $\mathbf{z}_1$  and  $\mathbf{z}_2$ ,

$$\mathcal{L}_{ms} = \max_G \left( \frac{d_{\mathbf{I}}(G(\mathbf{z}_1, \mathbf{I}), G(\mathbf{I}, \mathbf{z}_2))}{d_{\mathbf{z}}(\mathbf{z}_1, \mathbf{z}_2)} \right), \quad (6)$$

where  $d_*(\cdot)$  denotes the distance metric.

The regularization term can be easily applied to the proposed framework:

$$\mathcal{L}_{new} = \mathcal{L}_{ori} + \lambda_{ms} \mathcal{L}_{ms}, \quad (7)$$

where  $\mathcal{L}_{ori}$  denote the full objective.

### 3.5 Multi-Domain Image-to-Image Translation

In addition to the translation between two domains, we apply the proposed disentangle representation to the multi-domain setting. Different from two-domain I2I, multi-domain I2I aims to perform translation among multiple domains with a single generator  $G$ .

We illustrate the framework for multi-domain I2I in Figure 5. Given  $k$  domains  $\{N_i\}_{i=1 \sim k}$ , two images  $(x, y)$  and their one-hot domain code  $(z_x^d, z_y^d)$  are randomly sampled ( $x \in N_n, y \in N_m, Z^d \subset \mathbb{R}^k$ ). We encode the images onto a shared content space  $\mathcal{C}$ , and domain-specific attribute spaces  $\{\mathcal{A}_i\}_{i=1 \sim k}$ .

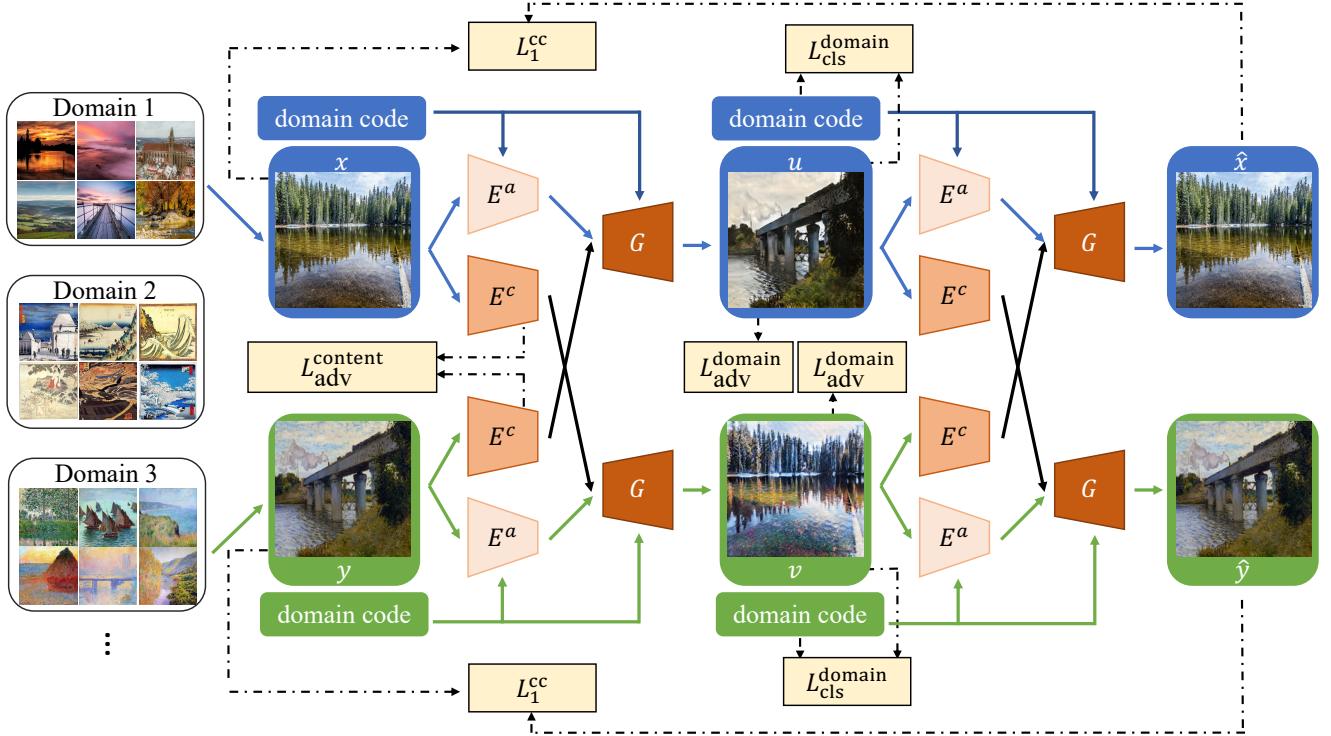
$$\begin{aligned} \{z_x^c, z_x^a\} &= \{E^c(x), E^a(x, z_x^d)\} & z_x^c \in \mathcal{C}, z_x^a \in \mathcal{A}_n \\ \{z_y^c, z_y^a\} &= \{E^c(y), E^a(y, z_y^d)\} & z_y^c \in \mathcal{C}, z_y^a \in \mathcal{A}_m \end{aligned} \quad (8)$$

We then perform the forward and backward translation similar to the two-domain translation.

$$\begin{aligned} u &= G(z_y^c, z_x^a, z_x^d) & v &= G(z_x^c, z_y^a, z_y^d) \\ \hat{x} &= G(z_v^c, z_u^a, z_u^d) & \hat{y} &= G(z_u^c, z_v^a, z_v^d) \end{aligned} \quad (9)$$

In addition to the loss functions used in the two-domain translation, we leverage the discriminator  $D$  as an auxiliary domain classifier. That is, the discriminator  $D$  not only aims to discriminate between real images and translated images ( $D_{dis}$ ), but also performs domain classification ( $D_{cls} : N_i \rightarrow Z^d$ ).

$$\begin{aligned} \mathcal{L}_{cls}^{domain} &= \mathbb{E}_{x, z_x^d} [-\log D_{cls}(z_x^d | x)] + \\ &\quad \mathbb{E}_{x, y, z_y^d} [-\log D_{cls}(z_y^d | G(z_x^c, z_y^a, z_y^d))] \end{aligned} \quad (10)$$



**Fig. 5: Multi-domains I2I framework.** We further extend the proposed disentangle representation framework to a more general multi-domain setting. Different from the class-specific encoders, generators, and discriminators used in two-domain I2I, all networks in mutli-domain are shared among all domains. Furthermore, one-hot domain codes are used as inputs and the discriminator will perform domain classification in addition to discrimination.

Therefore, our new objective functions are:

$$\begin{aligned} L_{D,D^c} = & \lambda_{\text{adv}}^{\text{content}} L_{\text{adv}}^c + \lambda_{\text{adv}}^{\text{domain}} L_{\text{adv}}^{\text{domain}} + \\ & \lambda_{\text{cls}}^{\text{domain}} \mathcal{L}_{\text{cls}}^{\text{domain}} \end{aligned} \quad (11)$$

$$\begin{aligned} L_{G,E^c,E^a} = & -L_{D,D^c} + \lambda_1^{\text{cc}} L_1^{\text{cc}} + \lambda_1^{\text{recon}} L_1^{\text{recon}} \\ & + \lambda_1^{\text{latent}} L_1^{\text{latent}} + \lambda_{\text{KL}} L_{\text{KL}} + \lambda_{\text{cls}}^{\text{domain}} \mathcal{L}_{\text{cls}}^{\text{domain}} \end{aligned} \quad (12)$$

## 4 Experimental Results

**Implementation details.** We implement our model with PyTorch [35]. We use the input image size of  $216 \times 216$  for all of our experiments. For the content encoder  $E^c$ , we use an architecture consisting of three convolution layers followed by four residual blocks. For the attribute encoder  $E^a$ , we use a CNN architecture with four convolution layers followed by fully-connected layers. We set the size of the attribute vector to  $z^a \in R^8$  for all experiments. For the generator  $G$ , we use an architecture consisting of four residual blocks followed by three fractionally strided convolution layers.

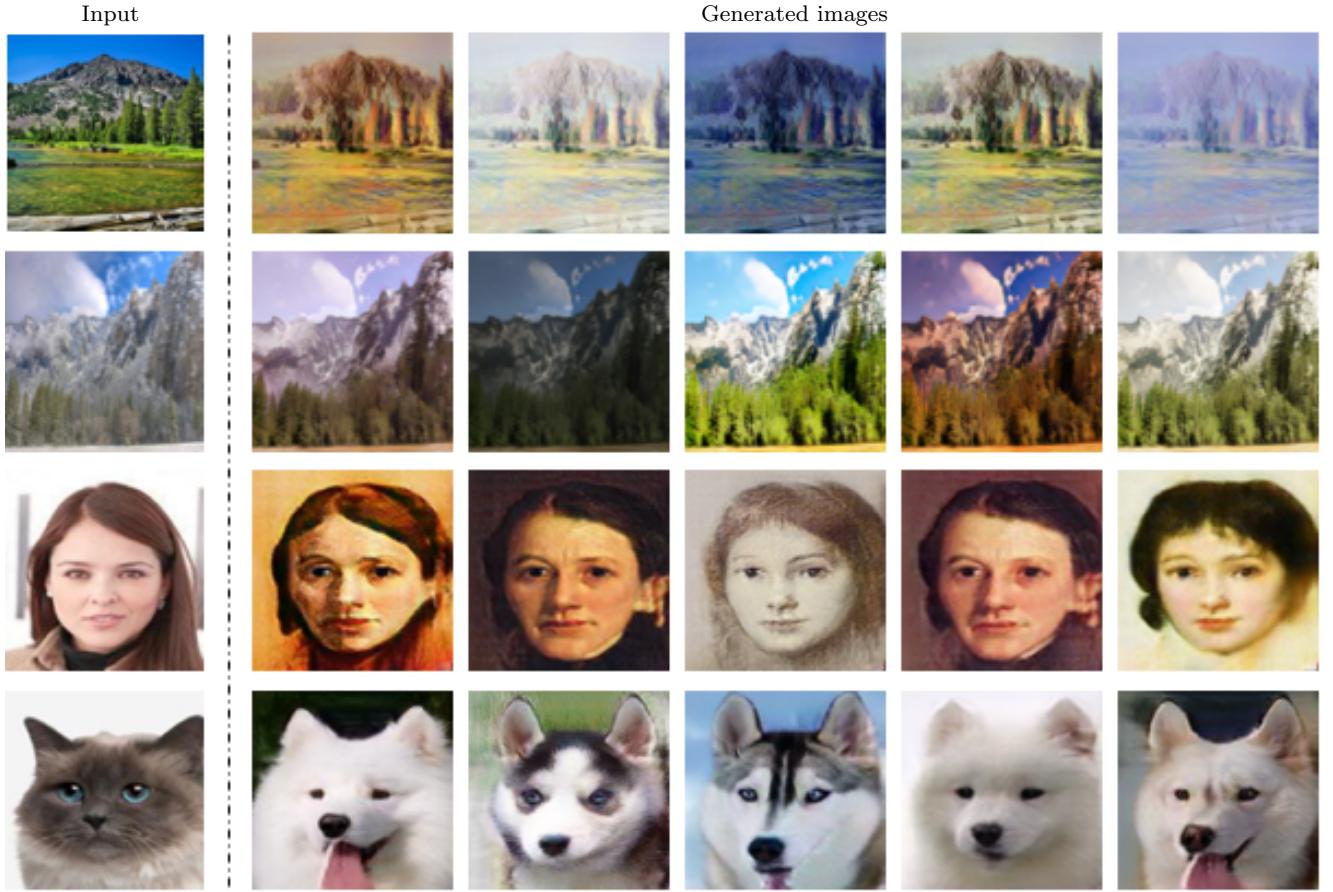
For training, we use the Adam optimizer [19] with a batch size of 1, a learning rate of 0.0001, and exponential decay rates  $(\beta_1, \beta_2) = (0.5, 0.999)$ . In all experiments, we set the hyper-parameters as follows:  $\lambda_{\text{adv}}^{\text{content}} = 1$ ,  $\lambda_{\text{adv}}^{\text{domain}} = 1$ ,  $\lambda_{\text{rec}}^{\text{content}} = 10$ ,  $\lambda_1^{\text{latent}} = 10$ , and  $\lambda_{\text{KL}} = 0.01$ . We also apply an L1 weight regularization on the content representation with a weight of 0.01. We follow the procedure in DCGAN [36] for training the model with adversarial loss.

**Datasets.** We evaluate our model on several datasets include Yosemite [49] (summer and winter scenes), pets (cat and dog) cropped from Google images, artworks [49] (Monet), and photo-to-portrait cropped from subsets of the WikiArt dataset <sup>1</sup> and the CelebA dataset [28].

**Compared methods.** We perform the evaluation on the following algorithms:

- **DRIT++:** The proposed model.
- **DRIT** [24], **MUNIT** [16]: Previous multimodal generation frameworks trained with unpaired data.
- **DRIT w/o  $D^c$ :** DRIT model without the content discriminator.
- **Cycle/Bicycle:** We construct a baseline using a combination of CycleGAN and BicycleGAN. Here,

<sup>1</sup> <https://www.wikiart.org/>



**Fig. 6: Sample results.** We show example results produced by our model. The left column shows the input images in the source domain. The other five columns show the output images generated by sampling random vectors in the attribute space. The mappings from top to bottom are: Photo  $\rightarrow$  Monet, winter  $\rightarrow$  summer, photograph  $\rightarrow$  portrait, and cat  $\rightarrow$  dog.

we first train CycleGAN on unpaired data to generate corresponding images as *pseudo* image pairs. We then use this pseudo paired data to train BicycleGAN.

– **CycleGAN** [49], **BicycleGAN** [50]

#### 4.1 Qualitative Evaluation

**Diversity.** We first demonstrate the visual artifacts of images generated by baseline methods in Figure 6. In Figure 7, we compare the proposed model with other methods. Both our model without  $D^c$  and Cycle/Bicycle can generate diverse results. However, the results contain clearly visible artifacts. Without the content discriminator, our model fails to capture domain-related details (e.g., the color of tree and sky). Therefore, the variations take place in global color difference. Cycle/Bicycle is trained on pseudo paired data generated by CycleGAN. The quality of the pseudo paired data is not uni-

formly ideal. As a result, the generated images are of ill-quality.

To have a better understanding of the learned domain-specific attribute space, we perform linear interpolation between two given attributes and generate the corresponding images as shown in Figure 9. The interpolation results validate the continuity in the attribute space and show that our model can generalize in the distribution, rather than memorize trivial visual information.

**Mode seeking regularization.** We demonstrate the effectiveness of the mode seeking regularization in Figure 8. The mode seeking regularization term substantially alleviate the mode collapse issue in DRIT [24], particularly in the challenging shape-variation translation (i.e., dog-to-cat translation).

**Attribute transfer.** We demonstrate the results of the attribute transfer in Figure 10. Thanks to the representation disentanglement of content and attribute,

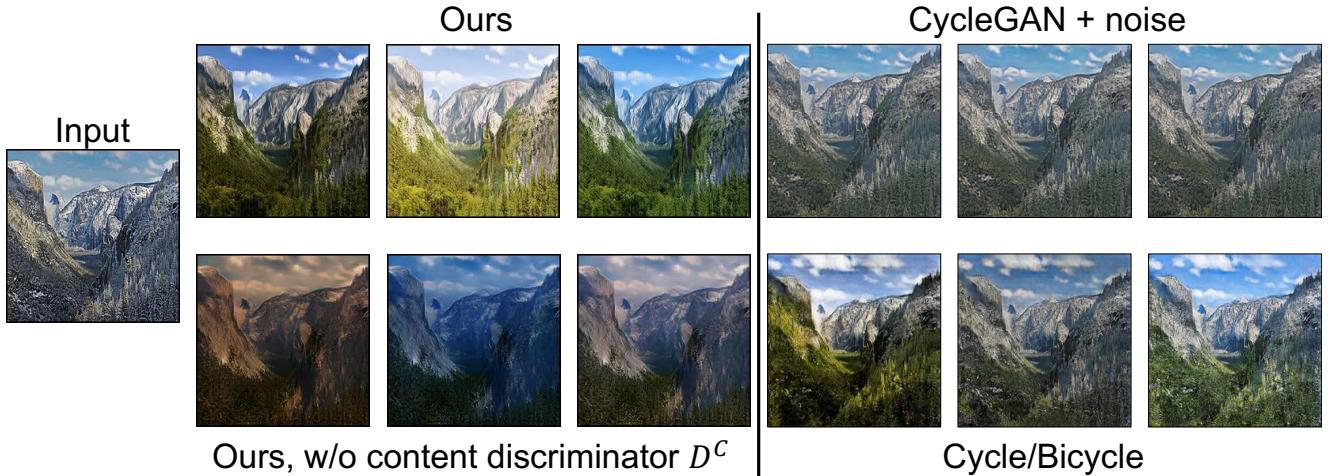


Fig. 7: **Baseline artifacts.** On the winter → summer translation task, our model produces more diverse and realistic samples over baselines.

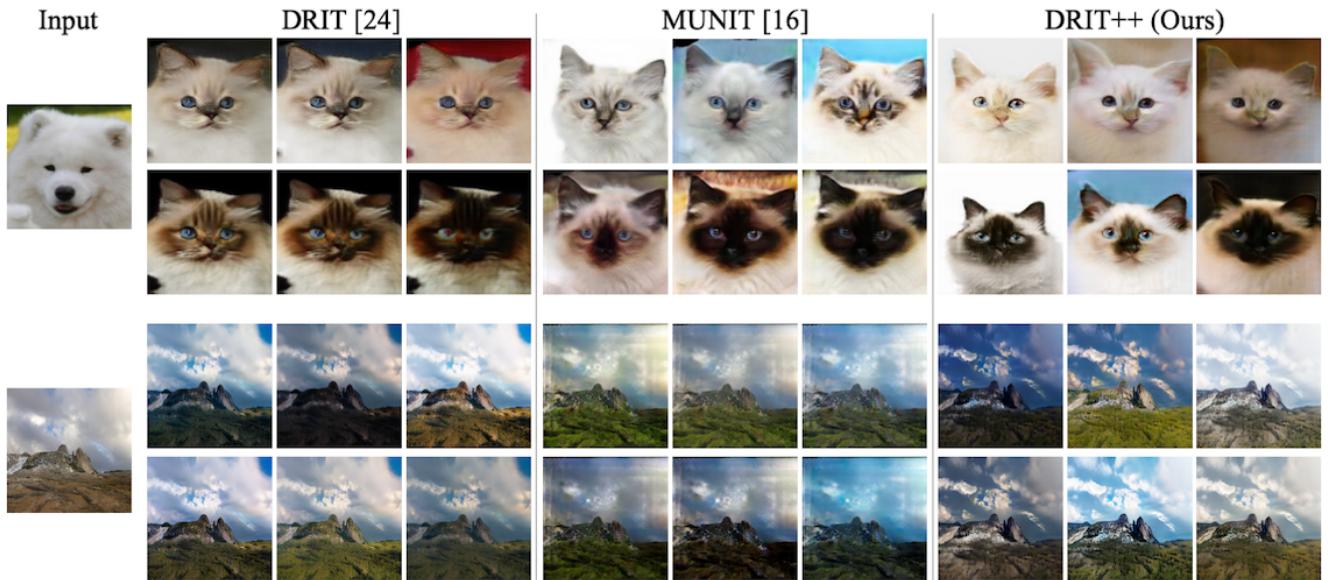


Fig. 8: **Effectiveness of mode seeking regularization.** Mode seeking regularization helps improve the diversity of translated images while maintaining the visual quality.

we are able to perform attribute transfer from images of desired attributes, as illustrated in Figure 3(c). Moreover, since the content space is shared between two domains, we can generate images conditioned on content features encoded from either domain. Thus our model can achieve not only inter-domain but also intra-domain attribute transfer. Note that intra-domain attribute transfer is not explicitly involved in the training process.

**Multi-domain I2I.** Figure 5 shows the results of applying the proposed method on the multi-domain I2I. We perform translation among three domains (real images and two artistic styles) and four domains (different weather conditions). Using a single generator, the

proposed model is able to perform diverse translation among multiple domains.

#### 4.2 Quantitative Evaluation

**Metrics** We conduct quantitative evaluations using the following metrics.

- **FID.** To evaluate the quality of the generated images, we use FID [14] to measure the distance between the generated distribution and the real one through features extracted by Inception Network [41]. Lower FID values indicate better quality of the generated images.
- **LPIPS.** To evaluate diversity, we employ LPIPS [48]. LPIPS measures the average feature distances be-

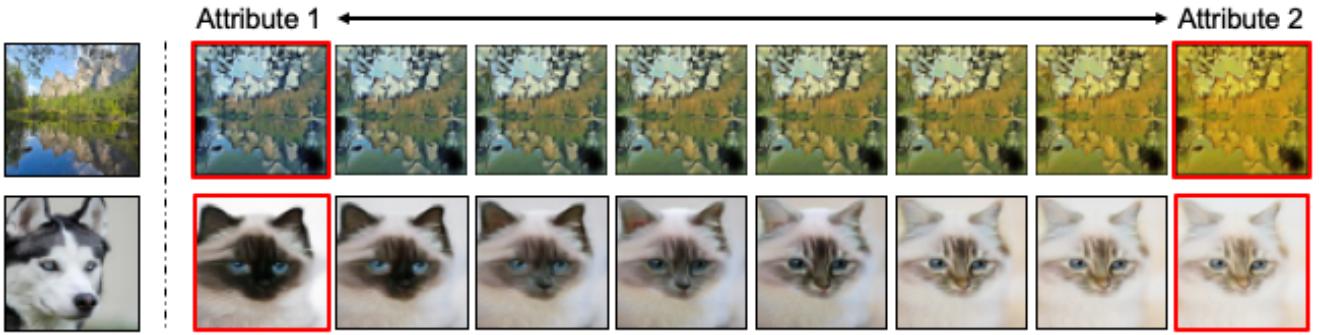


Fig. 9: **Linear interpolation between two attribute vectors.** Translation results with linear-interpolated attribute vectors between two attributes (highlighted in red).

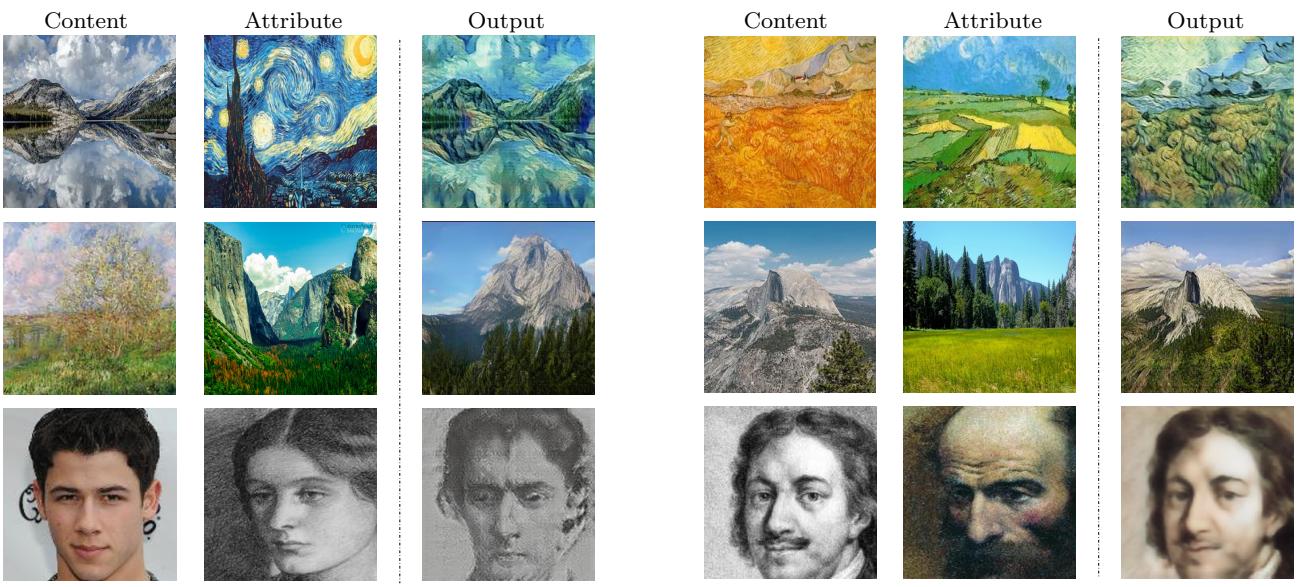


Fig. 10: **Attribute transfer.** At test time, in addition to random sampling from the attribute space, we can also perform translation with the query images with the desired attributes. Since the content space is shared across the two domains, we not only can achieve (a) inter-domain, but also (b) intra-domain attribute transfer. Note that we do not explicitly involve intra-domain attribute transfer during training.

tween generated samples. Higher LPIPS score indicates better diversity among the generated images.

- **JSD and NDB.** To measure the similarity between the distribution between real images and generated one, we adopt two bin-based metrics, JSD and NDB, proposed in [38]. These metrics evaluate the extent of mode missing of generative models. Following [38], we first cluster the training samples using K-means into different bins. These bins can be viewed as modes of the real data distribution. We then assign each generated sample to the bin of its nearest neighbor. We calculate the bin-proportions of the training samples and the synthesized samples to evaluate the difference between the generated dis-

tribution and the real data distribution. NDB score and JSD of the bin-proportion are then computed to measure the level of mode collapse. Lower NDB score and JSD mean the generated data distribution approaches the real data distribution better by fitting more modes. Please refer to [38] for more details.

- **User preference.** For evaluating realism, we conduct a user study using pairwise comparison. Given a pair of images sampled from real images and translated images generated from various methods, users need to answer the question “Which image is more realistic?”

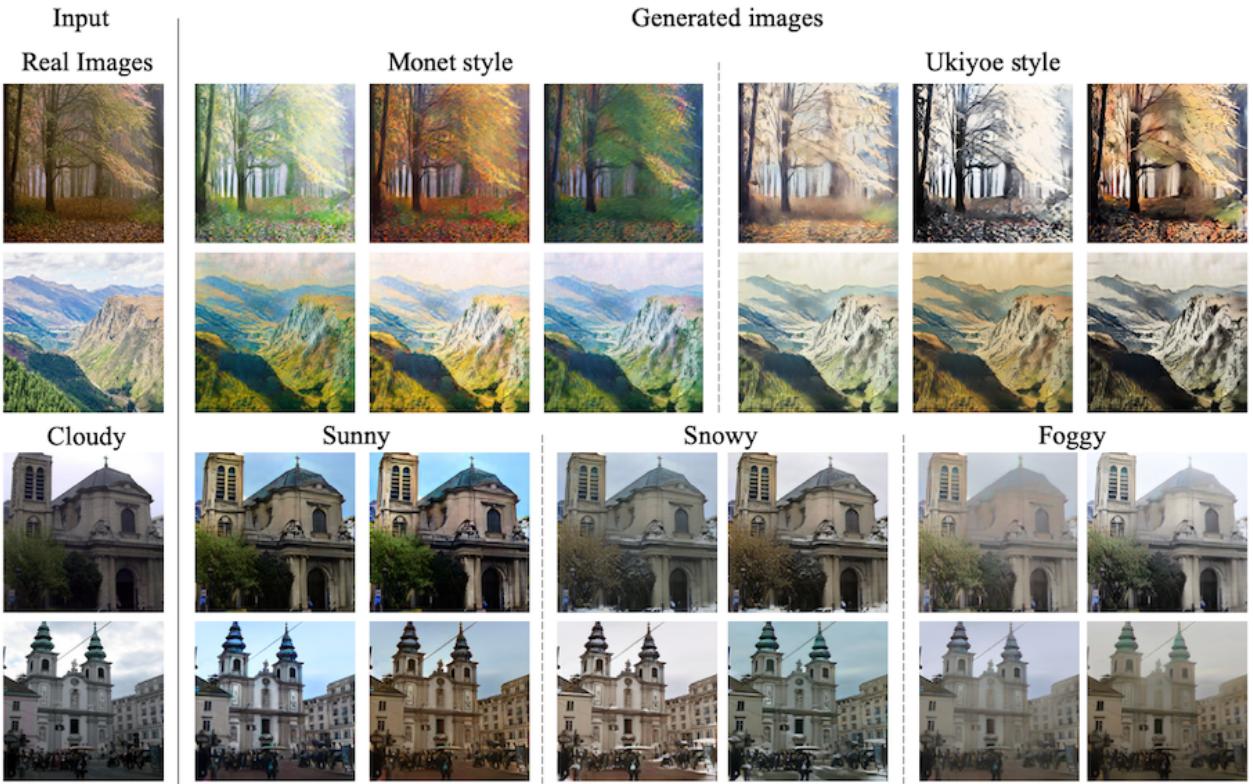


Fig. 11: **Multi-domain I2I.** We show example results of our model on the multi-domain I2I task. We demonstrate the translation among real images and two artistic styles (Monet and Ukiyoe), and the translation among different weathers (sunny, cloudy, snowy, and foggy).

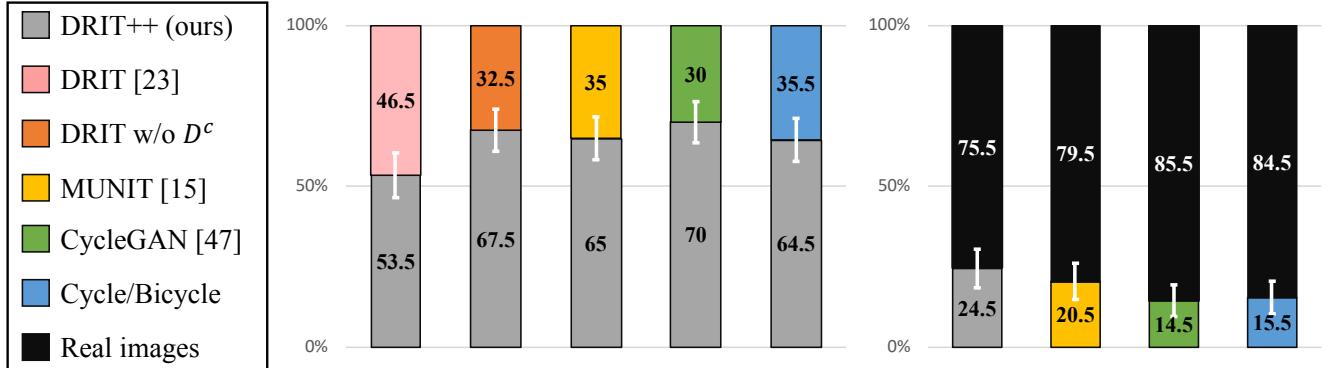


Fig. 12: **Realism preference results.** We conduct a user study to ask subjects to select results that are *more realistic* through pairwise comparisons. The number indicates the percentage of preference for that comparison pair. We use the winter → summer and the cat → dog translation for this experiment.

**Realism vs. diversity.** We conduct the experiment using winter → summer and cat → dog translation with the Yosemite and pets datasets, respectively. Table 1, Table 2, and Figure 12 present the quantitative comparisons with other methods as well as baseline methods. In Table 1, DRIT++ excels at all metrics. DRIT++ generates images that are not only realistic, but also diverse and close to the original data distribution. Table 2 validates the effectiveness of the content discriminator, the latent regression loss, and the mode-seeking reg-

ularization. Figure 12 shows the results of user study. DRIT++ outperforms previous work as well as baseline methods.

### 4.3 High Resolution I2I

We demonstrate that the proposed scheme can be applied to the translation tasks with high-resolution images. We perform the translation on the street scene (GTA [39] ↔ Cityscape [11]) dataset. The size of the

Table 1: Quantitative results of the Yosemite (Summer $\rightleftharpoons$ Winter) and the Cat $\rightleftharpoons$ Dog dataset.

Datasets		Winter $\rightarrow$ Summer			
		Cycle/Bicycle	DRIT	MUNIT	DRIT++
FID $\downarrow$		67.04 $\pm$ 0.60	41.34 $\pm$ 0.20	57.09 $\pm$ 0.37	<b>41.02 <math>\pm</math> 0.24</b>
NDB $\downarrow$		9.36 $\pm$ 0.69	9.38 $\pm$ 0.74	9.53 $\pm$ 0.64	<b>9.22 <math>\pm</math> 0.97</b>
JSD $\downarrow$		0.290 $\pm$ 0.086	0.304 $\pm$ 0.075	0.293 $\pm$ 0.062	<b>0.222 <math>\pm</math> 0.070</b>
LPIPS $\uparrow$		0.0974 $\pm$ 0.0003	0.0965 $\pm$ 0.0004	0.1136 $\pm$ 0.0008	<b>0.1183 <math>\pm</math> 0.0007</b>
Datasets		Cat $\rightarrow$ Dog			
		Cycle/Bicycle	DRIT	MUNIT	DRIT++
FID $\downarrow$		54.008 $\pm$ 1.590	24.306 $\pm$ 0.329	22.127 $\pm$ 0.712	<b>17.253 <math>\pm</math> 0.648</b>
NDB $\downarrow$		9.23 $\pm$ 0.84	8.16 $\pm$ 1.60	8.21 $\pm$ 1.17	<b>7.57 <math>\pm</math> 1.25</b>
JSD $\downarrow$		0.262 $\pm$ 0.072	0.075 $\pm$ 0.046	0.132 $\pm$ 0.066	<b>0.041 <math>\pm</math> 0.014</b>
LPIPS $\uparrow$		0.147 $\pm$ 0.001	0.245 $\pm$ 0.002	0.244 $\pm$ 0.002	<b>0.280 <math>\pm</math> 0.002</b>

Table 2: Ablation study.

	DRIT w/o $D^c$	DRIT w/o KL	DRIT w/o $L_1^{\text{latent}}$	DRIT	DRIT++
FID $\downarrow$	46.92 $\pm$ 0.35	<b>40.08 <math>\pm</math> 0.33</b>	53.12 $\pm$ 0.16	41.34 $\pm$ 0.20	41.02 $\pm$ 0.24
NDB $\downarrow$	9.36 $\pm$ 0.72	9.47 $\pm$ 0.70	9.97 $\pm$ 0.17	9.38 $\pm$ 0.74	<b>9.22 <math>\pm</math> 0.97</b>
JSD $\downarrow$	0.277 $\pm$ 0.077	0.289 $\pm$ 0.066	0.494 $\pm$ 0.045	0.304 $\pm$ 0.075	<b>0.222 <math>\pm</math> 0.070</b>
LPIPS $\uparrow$	0.0954 $\pm$ 0.0006	0.0957 $\pm$ 0.0007	0.0158 $\pm$ 0.0003	0.0965 $\pm$ 0.0004	<b>0.1183 <math>\pm</math> 0.0007</b>

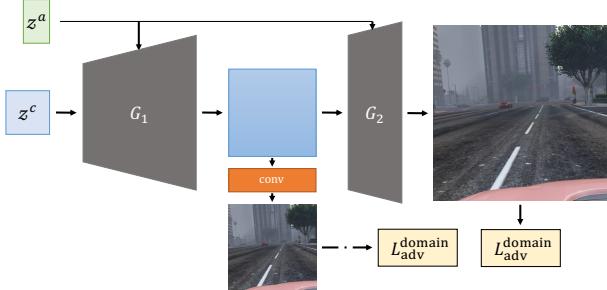


Fig. 13: **Multi-scale generator-discriminator.** To enhance the quality of generated high-resolution images, we adopt a multi-scale generator-discriminator architecture. We generate low-resolution images from the intermediate features of the generator. An additional adversarial domain loss is applied on the low-resolution images.

input image is  $720 \times 360$ . During the training, we randomly crop the image to the size of  $340 \times 340$  for memory efficiency consideration. To enhance the quality of the generated high-resolution images, we adopt a multi-scale generator-discriminator structure similar to StackGAN [46]. As shown in Figure 13, we extract the intermediate feature of the generator and pass through a convolutional layer to generate low-resolution images. We utilize an additional discriminator which takes low-resolution images as input. This discriminator enforces the first few layers of the generator to capture the distribution of low-level variations such as colors and image structures. We find such multi-scale generator-discriminator structure facilitate the training and yields more realistic images on high-resolution translation task. As the ex-

ample results shown in Figure 14, our proposed scheme with the multi-scale architecture is capable of generating diverse high-resolution images.

#### 4.4 Limitations

Our method has the following limitations. First, due to the limited amount of training data, the attribute space is not fully exploited. Our I2I translation fails when the sampled attribute vectors locate in under-sampled space, see Figure 15(a). Second, it remains difficult when the domain characteristics differ significantly. For example, Figure 15(b) shows a failure case on the human figure due to the lack of human-related portraits in Monet collections. Third, we use multiple encoders and decoders for the cross-cycle consistency during training, which requires large memory usage. The memory usage limits the application on high-resolution image-to-image translation.

#### 5 Conclusions

In this paper, we present a novel disentangled representation framework for diverse image-to-image translation with unpaired data. we propose to disentangle the latent space to a content space that encodes common information between domains, and a domain-specific attribute space that can model the diverse variations given the same content. We apply a content discriminator to facilitate the representation disentanglement. We propose a cross-cycle consistency loss for cyclic reconstruction to train in the absence of paired data.



Fig. 14: **High-resolution translations.** We show the example results produced by our model with multi-scale generator-discriminator architecture. The mappings from top to bottom are: GTA  $\rightarrow$  Cityscape, Cityscape  $\rightarrow$  GTA.



Fig. 15: **Failure Cases.** Typical cases: (a) Attribute space not fully exploited. (b) Distribution characteristic difference.

Qualitative and quantitative results show that the proposed model produces realistic and diverse images. We also apply the proposed method to domain adaptation and achieve competitive performance compared to the state-of-the-art methods.

## References

- Almahairi, A., Rajeswar, S., Sordoni, A., Bachman, P., Courville, A.: Augmented cyclegan: Learning many-to-many mappings from unpaired data. arXiv preprint arXiv:1802.10151 (2018) [3](#)
- Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN. In: ICML (2017) [3](#)
- Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised pixel-level domain adaptation with generative adversarial networks. In: CVPR (2017) [1](#)
- Cao, J., Katzir, O., Jiang, P., Lischinski, D., Cohen-Or, D., Tu, C., Li, Y.: Dida: Disentangled synthesis for domain adaptation. arXiv preprint arXiv:1805.08019 (2018) [3](#)
- Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. In: ICCV (2017) [1](#)
- Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., Abbeel, P.: InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In: NIPS (2016) [4](#)
- Chen, Y.C., Lin, Y.Y., Yang, M.H., Huang, J.B.: Crdoco: Pixel-level domain transfer with cross-domain consistency. In: CVPR (2019) [1](#)
- Cheung, B., Livezey, J.A., Bansal, A.K., Olshausen, B.A.: Discovering hidden factors of variation in deep networks. In: ICLR workshop (2015) [4](#)
- Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: CVPR, vol. 1711 (2018) [1](#)
- Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: CVPR (2018) [4](#)
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: CVPR (2016) [11](#)
- Denton, E.L., Birodkar, V.: Unsupervised learning of disentangled representations from video. In: NIPS (2017) [4](#)
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NIPS (2014) [3](#)
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. In: NIPS (2017) [3, 9](#)
- Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A.A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. arXiv preprint arXiv:1711.03213 (2017) [1](#)
- Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: ECCV (2018) [2, 3, 7](#)
- Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: CVPR (2017) [1, 2, 3](#)
- Kim, T., Cha, M., Kim, H., Lee, J., Kim, J.: Learning to discover cross-domain relations with generative adversarial networks. In: ICML (2017) [3](#)
- Kingma, D., Adam, J.B.: A method for stochastic optimization. In: ICLR (2015) [7](#)

20. Kingma, D.P., Rezende, D., Mohamed, S.J., Welling, M.: Semi-supervised learning with deep generative models. In: NIPS (2014) [4](#)
21. Lai, W.S., Huang, J.B., Ahuja, N., Yang, M.H.: Deep laplacian pyramid networks for fast and accurate super-resolution. In: CVPR (2017) [1](#)
22. Larsson, G., Maire, M., Shakhnarovich, G.: Learning representations for automatic colorization. In: ECCV (2016) [1](#)
23. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., Shi, W.: Photo-realistic single image super-resolution using a generative adversarial network. In: CVPR (2017) [1, 3](#)
24. Lee, H.Y., Tseng, H.Y., Huang, J.B., Singh, M.K., Yang, M.H.: Diverse image-to-image translation via disentangled representations. In: ECCV (2018) [2, 7, 8](#)
25. Li, Y., Huang, J.B., Ahuja, N., Yang, M.H.: Deep joint image filtering. In: ECCV (2016) [1](#)
26. Liu, A., Liu, Y.C., Wang, F.Y.C.: A unified feature disentangler for multi-domain image translation and manipulation. In: NIPS (2018) [4](#)
27. Liu, M.Y., Breuel, T., Kautz, J.: Unsupervised image-to-image translation networks. In: NIPS (2017) [1, 2, 3, 4](#)
28. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: ICCV (2015) [7](#)
29. Ma, L., Jia, X., Georgoulis, S., Tuytelaars, T., Van Gool, L.: Exemplar guided unsupervised image-to-image translation. arXiv preprint arXiv:1805.11145 (2018) [3](#)
30. Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., Frey, B.: Adversarial autoencoders. In: ICLR workshop (2016) [4](#)
31. Mao, Q., Lee, H.Y., Tseng, H.Y., Ma, S., Yang, M.H.: Mode seeking generative adversarial networks for diverse image synthesis. In: IEEE Conference on Computer Vision and Pattern Recognition (2019) [3, 6](#)
32. Mathieu, M., Zhao, J., Sprechmann, P., Ramesh, A., LeCun, Y.: Disentangling factors of variation in deep representation using adversarial training. In: NIPS (2016) [4](#)
33. Murez, Z., Kolouri, S., Kriegman, D., Ramamoorthi, R., Kim, K.: Image to image translation for domain adaptation. In: CVPR (2018) [1](#)
34. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization (2019) [1](#)
35. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: NIPS workshop (2017) [7](#)
36. Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. In: ICLR (2016) [3, 7](#)
37. Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. In: ICML (2016) [3](#)
38. Richardson, E., Weiss, Y.: On GANs and GMMs. In: NIPS (2018) [3, 10](#)
39. Richter, S.R., Vineet, V., Roth, S., Koltun, V.: Playing for data: Ground truth from computer games. In: ECCV (2016) [11](#)
40. Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R.: Learning from simulated and unsupervised images through adversarial training. In: CVPR (2017) [1](#)
41. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: CVPR (2015) [9](#)
42. Taigman, Y., Polyak, A., Wolf, L.: Unsupervised cross-domain image generation. In: ICLR (2017) [1](#)
43. Vondrick, C., Pirsiavash, H., Torralba, A.: Generating videos with scene dynamics. In: NIPS (2016) [3](#)
44. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. In: CVPR (2018) [1](#)
45. Yi, Z., Zhang, H.R., Tan, P., Gong, M.: Dualgan: Unsupervised dual learning for image-to-image translation. In: ICCV (2017) [1](#)
46. Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., Metaxas, D.: Stackgan++: Realistic image synthesis with stacked generative adversarial networks. TPAMI (2018) [12](#)
47. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: ECCV (2016) [1](#)
48. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep networks as a perceptual metric. In: CVPR (2018) [3, 9](#)
49. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017) [1, 2, 3, 4, 7, 8](#)
50. Zhu, J.Y., Zhang, R., Pathak, D., Darrell, T., Efros, A.A., Wang, O., Shechtman, E.: Toward multimodal image-to-image translation. In: NIPS (2017) [2, 3, 6, 8](#)