

StarGAN v2: Diverse Image Synthesis for Multiple Domains

Yunjey Choi* Youngjung Uh* Jaejun Yoo* Jung-Woo Ha

Clova AI Research, NAVER Corp.

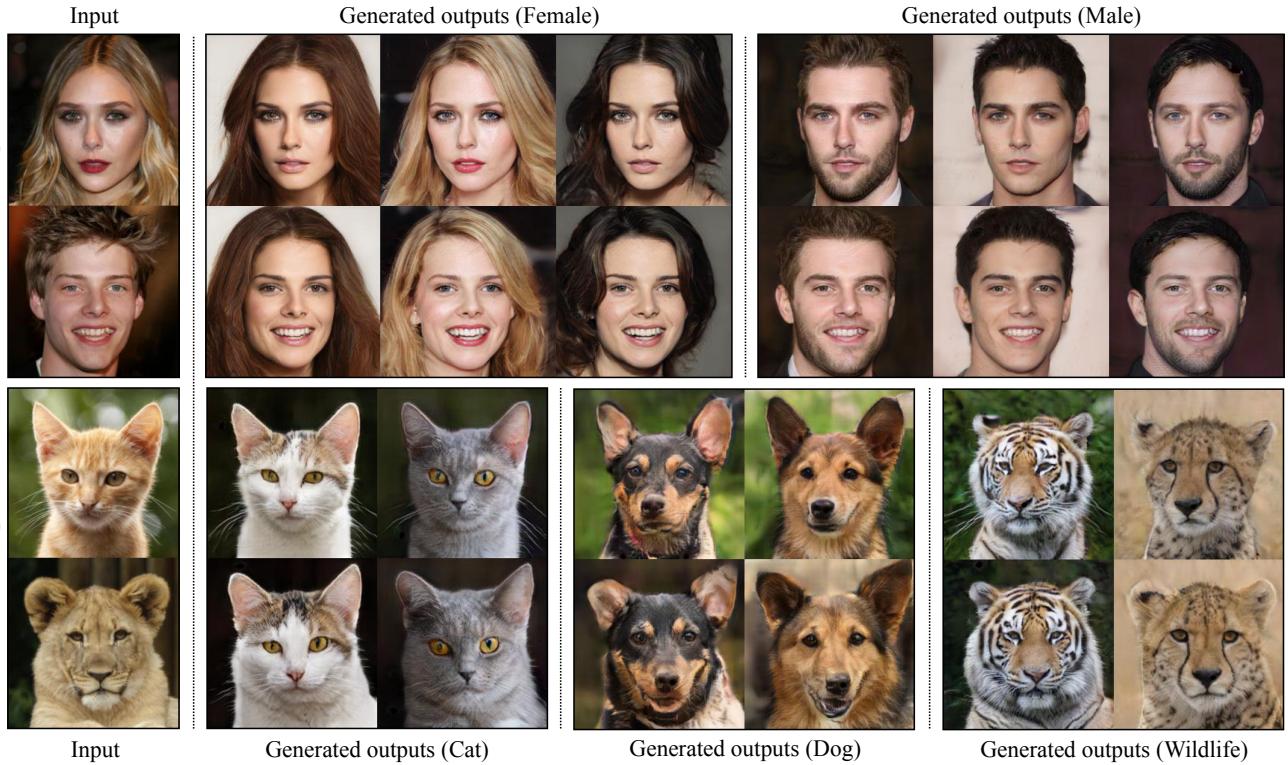


Figure 1. Diverse image synthesis results on the CelebA-HQ dataset and the newly collected animal faces (AFHQ) dataset. The first column shows input images while the remaining columns are images synthesized by StarGAN v2.

Abstract

A good image-to-image translation model should learn a mapping between different visual domains while satisfying the following properties: 1) diversity of generated images and 2) scalability over multiple domains. Existing methods address either of the issues, having limited diversity or multiple models for all domains. We propose StarGAN v2, a single framework that tackles both and shows significantly improved results over the baselines. Experiments on CelebA-HQ and a new animal faces dataset (AFHQ) validate our superiority in terms of visual quality, diversity, and scalability. To better assess image-to-image translation models, we release AFHQ, high-quality animal faces with large inter- and intra-domain differences. The code, pretrained models, and dataset can be found at clovaai/stargan-v2.

1. Introduction

Image-to-image translation aims to learn a mapping between different visual domains [20]. Here, domain implies a set of images that can be grouped as a visually distinctive category, and each image has a unique appearance, which we call style. For example, we can set image domains based on the gender of a person, in which case the style include makeup, beard, and hairstyle (top half of Figure 1). An ideal image-to-image translation method should be able to synthesize images considering the diverse styles in each domain. However, designing and learning such models become complicated as there can be arbitrarily large number of styles and domains in the dataset.

* indicates equal contribution

To address the style diversity, much work on image-to-image translation has been developed [1, 16, 33, 27, 37, 53]. These methods inject a low-dimensional latent code to the generator, which can be randomly sampled from the standard Gaussian distribution. Their domain-specific decoders interpret the latent codes as recipes for various styles when generating images. However, because these methods have only considered a mapping between two domains, they are not scalable to the increasing number of domains. For example, having N domains, these methods require to train $N(N-1)$ generators to handle translations between each and every domain, limiting their practical usage.

To address the scalability, several studies have proposed a unified framework [2, 7, 17, 29]. StarGAN [7] is one of the earliest models, which learns the mappings between all available domains using a single generator. The generator takes a domain label as an additional input, and learns to transform an image into the corresponding domain. However, StarGAN still learns a deterministic mapping per each domain, which does not capture the multi-modal nature of the data distribution. This limitation comes from the fact that each domain is indicated by a predetermined label. Note that the generator receives a fixed label (*e.g.*, one-hot vector) as input, and thus it inevitably produces the same output per each domain, given a source image.

To get the best of both worlds, we propose StarGAN v2, a scalable approach that can generate diverse images across multiple domains. In particular, we start from StarGAN and replace its domain label with our proposed domain-specific style code that can represent diverse styles of a specific domain. To this end, we introduce two modules, a mapping network and a style encoder. The mapping network learns to transform random Gaussian noise into a style code, while the encoder learns to extract the style code from a given reference image. Considering multiple domains, both modules have multiple output branches, each of which provides style codes for a specific domain. Finally, utilizing these style codes, our generator learns to successfully synthesize diverse images over multiple domains (Figure 1).

We first investigate the effect of individual components of StarGAN v2 and show that our model indeed benefits from using the style code (Section 3.1). We empirically demonstrate that our proposed method is scalable to multiple domains and gives significantly better results in terms of visual quality and diversity compared to the leading methods (Section 3.2). Last but not least, we present a new dataset of animal faces (AFHQ) with high quality and wide variations (Appendix A) to better evaluate the performance of image-to-image translation models on large inter- and intra-domain differences. We release this dataset publicly available for research community.

2. StarGAN v2

In this section, we describe our proposed framework and its training objective functions.

2.1. Proposed framework

Let \mathcal{X} and \mathcal{Y} be the sets of images and possible domains, respectively. Given an image $\mathbf{x} \in \mathcal{X}$ and an arbitrary domain $y \in \mathcal{Y}$, our goal is to train a single generator G that can generate diverse images of each domain y that corresponds to the image \mathbf{x} . We generate domain-specific style vectors in the learned style space of each domain and train G to reflect the style vectors. Figure 2 illustrates an overview of our framework, which consists of four modules described below.

Generator (Figure 2a). Our generator G translates an input image \mathbf{x} into an output image $G(\mathbf{x}, \mathbf{s})$ reflecting a domain-specific style code \mathbf{s} , which is provided either by the mapping network F or by the style encoder E . We use adaptive instance normalization (Adain) [15, 22] to inject \mathbf{s} into G . We note that \mathbf{s} is designed to represent a style of a specific domain y , which removes the necessity of providing y to G and allows G to synthesize images of all domains.

Mapping network (Figure 2b). Given a latent code \mathbf{z} and a domain y , our mapping network F generates a style code $\mathbf{s} = F_y(\mathbf{z})$, where $F_y(\cdot)$ denotes an output of F corresponding to the domain y . F consists of an MLP with multiple output branches to provide style codes for all available domains. F can produce diverse style codes by sampling the latent vector $\mathbf{z} \in \mathcal{Z}$ and the domain $y \in \mathcal{Y}$ randomly. Our multi-task architecture allows F to efficiently and effectively learn style representations of all domains.

Style encoder (Figure 2c). Given an image \mathbf{x} and its corresponding domain y , our encoder E extracts the style code $\mathbf{s} = E_y(\mathbf{x})$ of \mathbf{x} . Here, $E_y(\cdot)$ denotes the output of E corresponding to the domain y . Similar to F , our style encoder E benefits from the multi-task learning setup. E can produce diverse style codes using different reference images. This allows G to synthesize an output image reflecting the style \mathbf{s} of a reference image \mathbf{x} .

Discriminator (Figure 2d). Our discriminator D is a multi-task discriminator [29, 34], which consists of multiple output branches. Each branch D_y learns a binary classification determining whether an image \mathbf{x} is a real image of its domain y or a fake image $G(\mathbf{x}, \mathbf{s})$ produced by G .

2.2. Training objectives

Given an image $\mathbf{x} \in \mathcal{X}$ and its original domain $y \in \mathcal{Y}$, we train our framework using the following objectives.

Adversarial objective. During training, we sample a latent code $\mathbf{z} \in \mathcal{Z}$ and a target domain $\tilde{y} \in \mathcal{Y}$ randomly, and

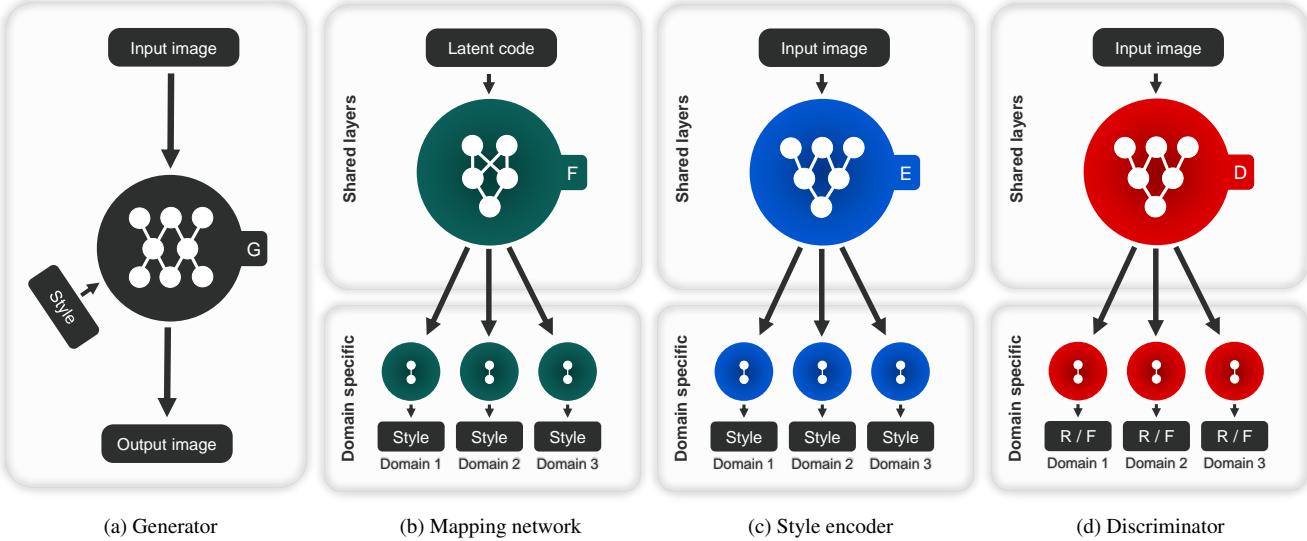


Figure 2. Overview of StarGAN v2, consisting of four modules. (a) The generator translates an input image into an output image reflecting the domain-specific style code. (b) The mapping network transforms a latent code into style codes for multiple domains, one of which is randomly selected during training. (c) The style encoder extracts the style code of an image, allowing the generator to perform reference-guided image synthesis. (d) The discriminator distinguishes between real and fake images from multiple domains.

generate a target style code $\tilde{s} = F_{\tilde{y}}(\mathbf{z})$. The generator G takes an image \mathbf{x} and \tilde{s} as inputs and learns to generate an output image $G(\mathbf{x}, \tilde{s})$ via adversarial loss

$$\mathcal{L}_{adv} = \mathbb{E}_{\mathbf{x}, y} [\log D_y(\mathbf{x})] + \mathbb{E}_{\mathbf{x}, \tilde{y}, \mathbf{z}} [\log (1 - D_{\tilde{y}}(G(\mathbf{x}, \tilde{s})))], \quad (1)$$

where $D_y(\cdot)$ denotes the output of D corresponding to the domain y . The mapping network F learns to provide the style code \tilde{s} that is likely in the target domain \tilde{y} , and G learns to utilize \tilde{s} and generate an image $G(\mathbf{x}, \tilde{s})$ that is indistinguishable from real images of the domain \tilde{y} .

Style reconstruction. In order to enforce the generator G to utilize the style code \tilde{s} when generating the image $G(\mathbf{x}, \tilde{s})$, we employ a style reconstruction loss

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

This objective is similar to the previous approaches [16, 27], which employ multiple encoders to learn a mapping from an image to its latent code. The notable difference is that we train a single encoder E to encourage diverse outputs for multiple domains. At test time, our learned encoder E allows G to transform an input image, reflecting the style of a reference image.

Style diversification. To further enable the generator G to produce diverse images, we explicitly regularize G with the diversity sensitive loss [33, 47]

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

where the target style codes \tilde{s}_1 and \tilde{s}_2 are produced by F conditioned on two random latent codes \mathbf{z}_1 and \mathbf{z}_2 (*i.e.*,

$\tilde{s}_i = F_{\tilde{y}}(\mathbf{z}_i)$ for $i \in \{1, 2\}$). Maximizing the regularization term forces G to explore the image space and discover meaningful style features to generate diverse images. We adopt the ℓ_1 -norm as our distance metric. Since the objective does not have an optimal point, we linearly decay the weight of the loss to zero during training.

Preserving source characteristics. The objectives (Eqs. (1) ~ (3)) do not guarantee that the generated image $G(\mathbf{x}, \tilde{s})$ properly preserves domain-invariant characteristics (*e.g.*, pose) of its input image \mathbf{x} . We address this issue by employing the cycle consistency loss [7, 23, 52]

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

where $\hat{s} = E_y(\mathbf{x})$ is the estimated style code of the input image \mathbf{x} , and y is the original domain of \mathbf{x} . By encouraging the generator G to reconstruct the input image \mathbf{x} with the estimated style code \hat{s} , G learns to preserve the original characteristics of \mathbf{x} while changing its style faithfully.

Full objective. Our full objective functions can be summarized as follows:

$$\mathcal{L}_D = -\mathcal{L}_{adv}, \quad (5)$$

$$\begin{aligned} \mathcal{L}_{F,G,E} = & \mathcal{L}_{adv} + \lambda_{sty} \mathcal{L}_{sty} \\ & - \lambda_{ds} \mathcal{L}_{ds} + \lambda_{cyc} \mathcal{L}_{cyc}, \end{aligned} \quad (6)$$

where λ_{sty} , λ_{ds} , and λ_{cyc} are hyperparameters that control the relative importance of each loss compared to the adversarial loss. We provide the detailed settings in Appendix B.

Method	FID	LPIPS
A Baseline StarGAN [7]	98.4	-
B + Multi-task discriminator	91.4	-
C + Tuning (e.g., R_1 regularization)	80.5	-
D + Latent code injection	32.3	0.312
E + Replace (D) with style code	21.2	0.406
F + Diversity regularization	18.0	0.428

Table 1. Performance of various configurations on CelebA-HQ. Frechet inception distance (FID) indicates the distance between two distributions of real and generated images (lower is better), while learned perceptual image patch similarity (LPIPS) measures the diversity of generated images (higher is better).

3. Experiments

In this section, we describe evaluation setups and conduct a set of experiments. We analyze the individual components of StarGAN v2 (Section 3.1) and compare our model with three leading baselines on diverse image synthesis (Section 3.2). All experiments are conducted using unseen images during the training phase.

Baselines. We use MUNIT [16], DRIT [27], and MSGAN [33] as our baselines, all of which learn multi-modal mappings between two domains. For multi-domain comparisons, we train these models multiple times for every pair of image domains. We also compare our method with StarGAN [7], which learns mappings among multiple domains using a single generator. All the baselines are trained using the implementations provided by the authors.

Datasets. We evaluate StarGAN v2 on CelebA-HQ [21] and our new AFHQ dataset (Appendix A). We separate CelebA-HQ into two domains of male and female, and AFHQ into three domains of cat, dog, and wildlife. Other than the domain labels, we do not use any additional information (e.g., facial attributes of CelebA-HQ or breeds of AFHQ) and let the models learn such information as styles without supervision. For a fair comparison, all images are resized to 256×256 resolution for training, which is the highest resolution used in the baselines.

Evaluation metrics. We evaluate both the visual quality and the diversity of generated images using Frechét inception distance (FID) [14] and learned perceptual image patch similarity (LPIPS) [51]. We compute FID and LPIPS for every pair of image domains within a dataset and report their average values. The details on evaluation metrics and protocols are further described in Appendix C.

3.1. Analysis of individual components

We evaluate individual components that are added to our baseline StarGAN using CelebA-HQ. Table 1 gives FID and LPIPS for several configurations, where each component is cumulatively added on top of StarGAN. An input image and the corresponding generated images of each con-

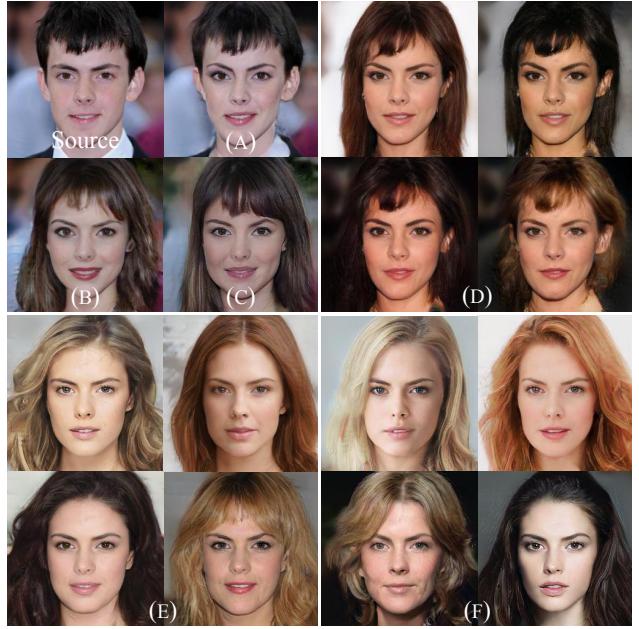


Figure 3. Visual comparison of generated images using each configuration in Table 1. Note that given a source image, the configurations (A) - (C) provide a single output, while (D) - (F) generate multiple output images.

figuration are shown in Figure 3. The baseline configuration (A) corresponds to the basic setup of StarGAN, which employs WGAN-GP [11], ACGAN discriminator [38], and depth-wise concatenation [35] for providing the target domain information to the generator. As shown in Figure 3a, the original StarGAN produces only a local change by applying makeup on the input image.

We first improve our baseline by replacing the ACGAN discriminator with a multi-task discriminator [34, 29], which allows the generator to transform the global structure of an input image as shown in configuration (B). Exploiting the recent advances in GANs, we further enhance the training stability and construct a new baseline (C) by applying R_1 regularization [34] and switching the depth-wise concatenation to adaptive instance normalization (AdaIN) [9, 15]. Note that we do not report LPIPS of these variations in Table 1, since they are yet to be designed to produce multiple outputs for a given input image and a target domain.

To induce diversity, one can easily think of giving a latent code z into the generator, where z is sampled from the standard Gaussian distribution [16, 27, 42, 46]. However, in a multi-domain scenario, we observe that this baseline (D) does not encourage the network to learn meaningful styles and fails to provide as much diversity as we expect. The model usually learns to modify intensity of pixels in a homogeneous region and still finds it difficult to change the overall structure (e.g., hairstyle). To address this, we introduce a mapping network with multiple output branches,

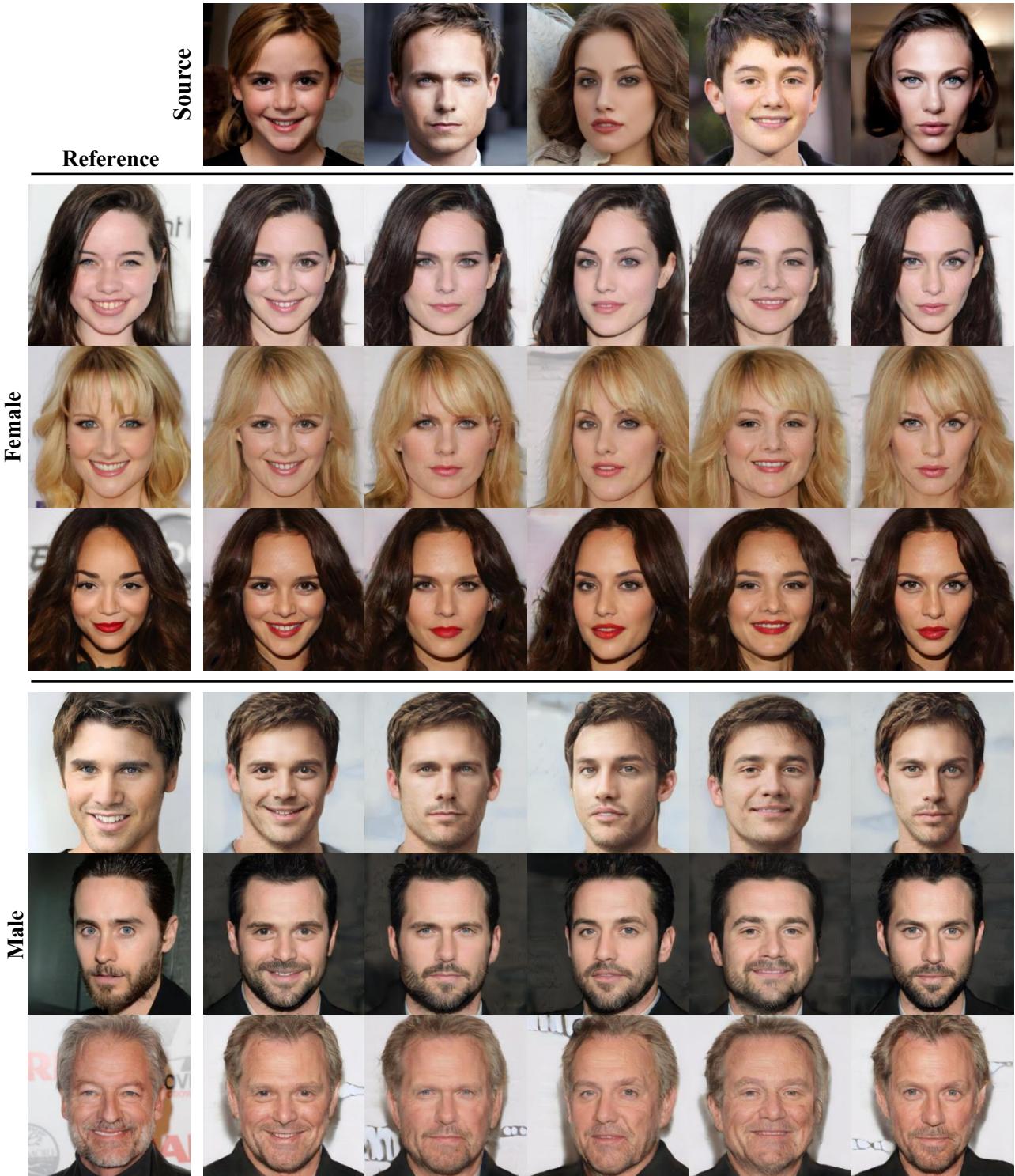
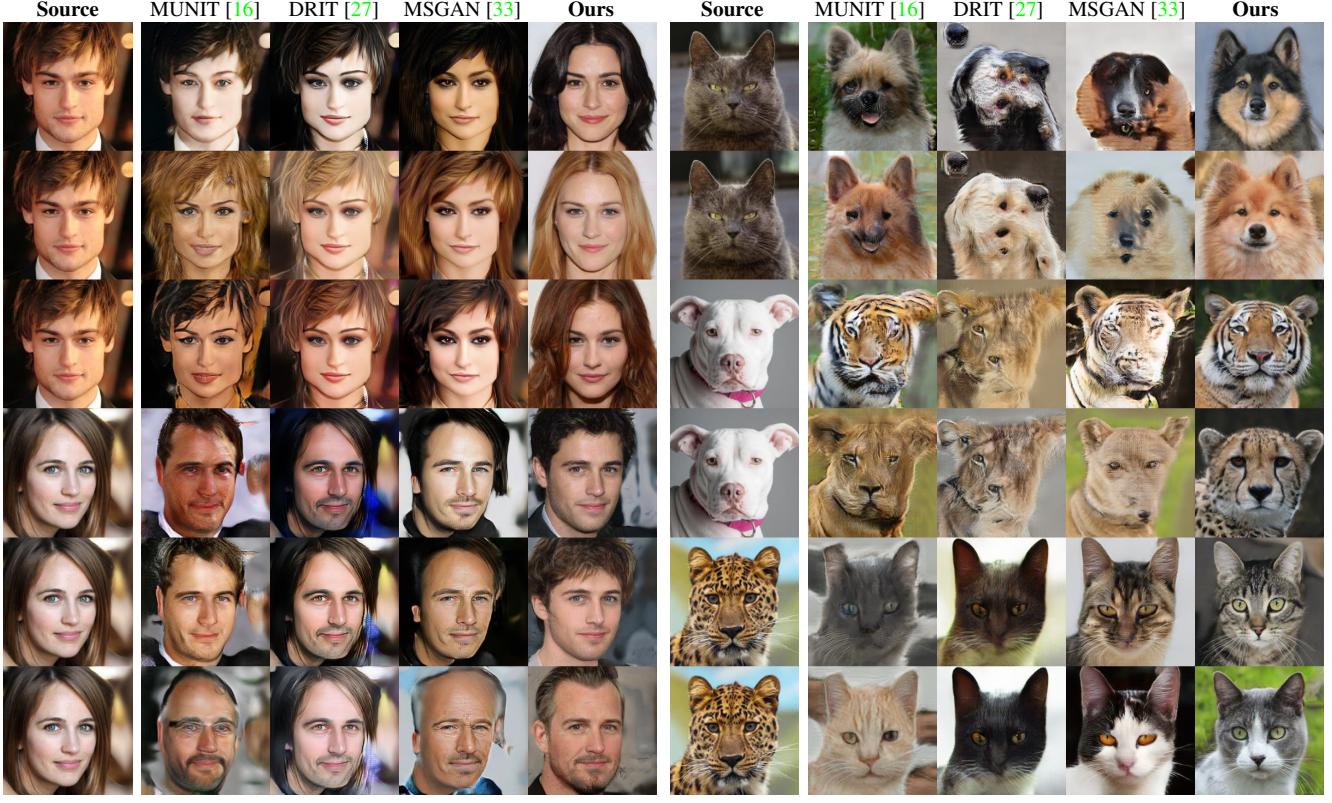


Figure 4. Reference-guided image synthesis results on CelebA-HQ. The source and reference images in the first row and the first column are real images, while the rest are images generated by our proposed model, StarGAN v2. Our model learns to transform a source image reflecting the style of a given reference image. High-level semantics such as hairstyle, makeup, beard and age are followed from the reference images, while the pose and identity of the source images are preserved. Note that the images in each column share a single identity with different styles, and those in each row share a style with different identities.



(a) Latent-guided synthesis on CelebA-HQ

(b) Latent-guided synthesis on AFHQ

Figure 5. Qualitative comparison of latent-guided image synthesis results on the [CelebA-HQ](#) and [AFHQ](#) datasets. Each method translates the source images (left-most column) to target domains using randomly sampled latent codes. (a) The top three rows correspond to the results of converting male to female and vice versa in the bottom three rows. (b) Every two rows from the top show the synthesized images in the following order: cat-to-dog, dog-to-wildlife, and wildlife-to-cat.

each of which provides style codes for a specific domain (E). We now explicitly have different style codes per domain that are injected to the generator (*i.e.*, $G(\mathbf{x}, \mathbf{s})$ instead of $G(\mathbf{x}, \mathbf{z}, \mathbf{y})$). This allows the generator to focus on utilizing a ready-made domain-specific style code instead of learning to combine a latent code with a domain label. Finally, we further improve the network to produce diverse outputs by adopting the diversity regularization (F), and this configuration corresponds to our proposed method, StarGAN v2. As shown in Figure 4, our method can synthesize images that reflect diverse styles of given reference images including hairstyle, makeup, beard, and age, without hurting the source characteristics. Additional results are presented in Appendix D.

3.2. Comparison on diverse image synthesis

In this section, we evaluate StarGAN v2 on diverse image synthesis from two perspectives: [latent-guided synthesis](#) and [reference-guided synthesis](#).

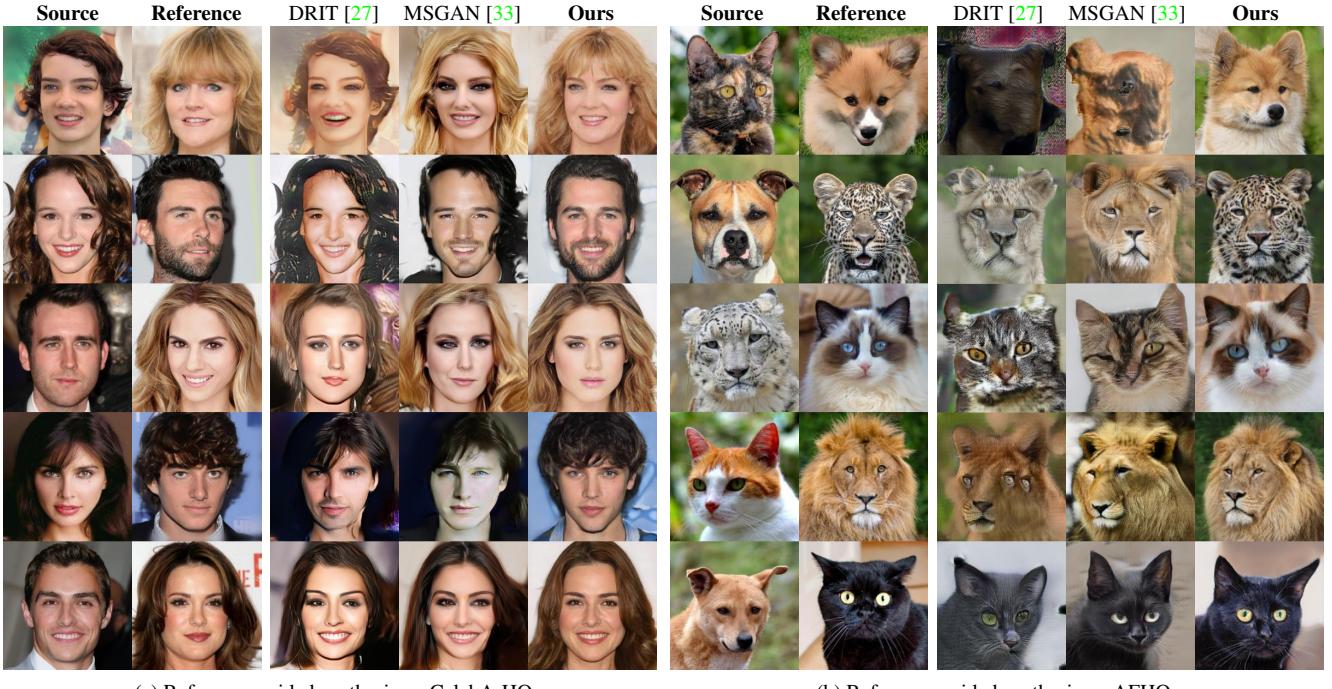
Latent-guided synthesis. Figure 5 provides a qualitative comparison of the competing methods. Each method produces multiple outputs using random Gaussian noises. For

Method	CelebA-HQ		AFHQ	
	FID	LPIPS	FID	LPIPS
MUNIT [16]	31.4	0.363	41.5	0.511
DRIT [27]	52.1	0.178	95.6	0.326
MSGAN [33]	33.1	0.389	61.4	0.517
StarGAN v2	18.0	0.428	24.4	0.524
Real images	15.1	-	13.1	-

Table 2. Quantitative comparison on latent-guided synthesis. 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.

CelebA-HQ, we observe that our method synthesizes images with a higher visual quality compared to the baseline models. In addition, our method is the only model that can successfully change the entire hair styles of the source images, which requires non-trivial effort (*e.g.*, generating ears). For AFHQ, which has relatively large variations, the performance of the baselines is considerably degraded, while our method still produces images with high quality and diverse styles.

As shown in Table 2, our method outperforms all the baselines by a large margin in terms of visual quality. For both CelebA-HQ and AFHQ, our method achieves FID of



(a) Reference-guided synthesis on CelebA-HQ

(b) Reference-guided synthesis on AFHQ

Figure 6. Qualitative comparison of reference-guided image synthesis results on the CelebA-HQ and AFHQ datasets. Each method translates the source images into target domains, reflecting the styles of the reference images.

18.0 and 24.4, respectively, which are about two times improvements to the previous leading method. Our LPIPS is also the highest among the baseline models, which implies our model produces the most diverse results given a single input image. We conjecture that the small gaps in LPIPS are due to the various artifacts from the baseline models.

Reference-guided synthesis. To obtain the style code from a reference image, we sample test images from a target domain and feed them to the encoder network of each method. For CelebA-HQ (Figure 6a), our method successfully renders distinctive styles (*e.g.*, bangs, beard, makeup, and hairstyle), while the others mostly match the color distribution of reference images. For the more challenging AFHQ dataset (Figure 6b), the baseline models suffer from a large domain shift. They hardly reflect the style of each reference image and only match the domain. In contrast, our model renders distinctive styles (*e.g.*, breeds) of each reference image as well as its fur pattern and eye color. Note that StarGAN v2 produces high quality images across all domains and these results are from a single generator. Since the other baselines are trained individually for each pair of domains, the output quality fluctuates across domains. For example, in AFHQ (Figure 6b), the baseline models work reasonably well in dog-to-wildlife (2nd row) while they fail in cat-to-dog (1st row).

Table 3 shows FID and LPIPS of each method for reference guided synthesis. For both datasets, our method achieves FID of 20.2, and 19.7, which are about 2 \times and

Method	CelebA-HQ		AFHQ	
	FID	LPIPS	FID	LPIPS
MUNIT [16]	107.1	0.176	223.9	0.199
DRIT [27]	53.3	0.311	114.8	0.156
MSGAN [33]	39.6	0.312	69.8	0.375
StarGAN v2	20.2	0.397	19.7	0.503
Real images	15.1	-	13.1	-

Table 3. Quantitative comparison on reference-guided synthesis. We sample ten reference images to synthesize diverse images.

3.5 \times better than the previous leading method, respectively. The LPIPS of StarGAN v2 is also the highest among the competitors, which implies that our model produces the most diverse results considering the styles of reference images. Here, MUNIT and DRIT suffer from mode-collapse in AFHQ, which results in lower LPIPS and higher FID than other methods.

Human evaluation. We use the Amazon Mechanical Turk (AMT) to compare the user preferences of our method with baseline approaches. Given a pair of source and reference images, the AMT workers are instructed to select one among four image candidates from the methods, whose order is randomly shuffled. We ask separately which model offers the best image quality and which model best stylizes the input image considering the reference image. For each comparison, we randomly generate 100 questions, and each question is answered by 10 workers. We also ask each worker a few simple questions to detect unworthy workers. The number of total valid workers is 76. As shown in

Method	CelebA-HQ		AFHQ	
	Quality	Style	Quality	Style
MUNIT [16]	6.2	7.4	1.6	0.2
DRIT [27]	11.4	7.6	4.1	2.8
MSGAN [33]	13.5	10.1	6.2	4.9
StarGAN v2	68.9	74.8	88.1	92.1

Table 4. Votes from AMT workers for the most preferred method regarding visual quality and style reflection (%). StarGAN v2 outperforms the baselines with remarkable margins in all aspects.

Table 4, our method obtains the majority of votes in all instances, especially in the challenging AFHQ dataset and the question about style reflection. These results show that StarGAN v2 better extracts and renders the styles onto the input image than the other baselines.

4. Discussion

We suggest several reasons why StarGAN v2 can successfully synthesize images of diverse styles over multiple domains. First, following the insight of StyleGAN [22], our style space is produced by non-linear transformations from a Gaussian distribution. This provides more flexibility to our model compared to assuming a fixed prior distribution [16, 27, 33] as shown in Section 3.2. Second, our style code is separately generated per each domain by the multi-branch encoder and mapping networks. By doing so, our generator can only focus on using the style code, whose domain-specific information is already taken care of by the encoder or the mapping network as shown in Section 3.1. Last but not least, our modules benefit from fully exploiting training data from multiple domains. By design, the shared part of each module should learn domain-invariant features which induces the regularization effect, encouraging better generalization to unseen samples. To show that our model generalizes over the unseen images, we test a few samples from FFHQ [22] with our model trained on CelebA-HQ (Figure 7). Here, StarGAN v2 successfully captures styles of references and renders these styles correctly to the source images, all of which are from the FFHQ dataset.

5. Related work

Generative adversarial networks (GANs) [10] have shown impressive results in many computer vision tasks such as image synthesis [4, 30, 8], colorization [18, 49] and super-resolution [26, 45]. Along with improving the visual quality of generated images, their diversity also has been considered as an important objective which has been tackled by either devoted loss functions [33, 34] or architectural design [4, 22]. StyleGAN [22] introduces a mapping network that embeds an input latent code into an intermediate style space to better represent the factors of variation. However, these methods cannot manipulate real images directly since they tackle unconditional generation from random noise.

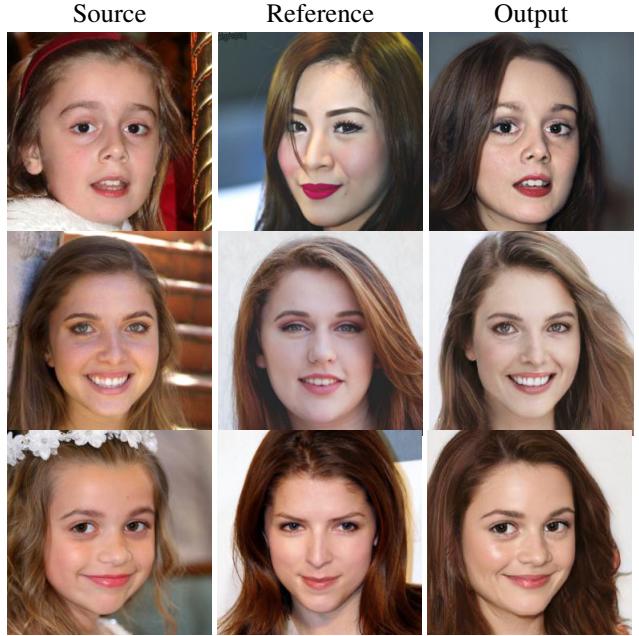


Figure 7. Reference-guided synthesis results on FFHQ with the model trained on CelebA-HQ. Despite the distribution gap between the two datasets, StarGAN v2 successfully extracts the style codes of the references and synthesizes faithful images.

Early methods for image-to-image translation are well known to learn a deterministic mapping even with stochastic noise inputs [20, 52, 28]. Several methods implicated a connection between stochastic noise and the output image for diversity, by either latent regression [53, 16, 27] or marginal matching [1]. Other approaches produce various outputs with the guidance of reference images [5, 6, 31, 39]. However, all these methods consider only two domains, and their extension to multiple domains is non-trivial. Recently, FUNIT [29] tackles multi-domain image translation using a few reference images from a target domain, but it does not offer latent-guided synthesis. Our method provides both latent-guided and reference-guided synthesis and is scalable to multiple domains. A concurrent method [50] tackles the same issue but they define the style as domain-shared characteristics, which limits the output diversity.

6. Conclusion

We address two major challenges in image-to-image translation; translating an image of one domain to diverse images of a target domain, and supporting multiple target domains. We proposed StarGAN v2 that tackles both issues within a single framework. The experimental results showed that our model can generate rich styles across different domains either by sampling a random latent vector or by using a reference image as a guide, remarkably outperforming the previous leading methods [16, 27, 33] in terms of visual quality and diversity. For evaluation on multiple domains, we gathered a new dataset of animal faces (AFHQ) that has large inter- and intra-domain variation.



Figure 8. Examples from our newly collected AFHQ dataset.

A. The AFHQ dataset

We release a new dataset of animal faces, **Animal Faces-HQ (AFHQ)**, consisting of 15,000 high-quality images at 512×512 resolution. Figure 8 shows example images of the AFHQ dataset. The dataset includes three domains of cat, dog, and wildlife, each providing 5000 images. By having multiple (three) domains and diverse images of various breeds (\geq eight) per each domain, AFHQ sets a more challenging image-to-image translation problem. For each domain, we select 500 images as a test set and provide all remaining images as a training set. We collected images with permissive licenses from the Flickr¹ and Pixabay² websites. All images are vertically and horizontally aligned to have the eyes at the center. The low-quality images were discarded by human effort. We will make the dataset available at <https://github.com/clovaai/stargan-v2>.

B. Training details

For fast training, the batch size is set to eight and the model is trained for 100K iterations. The training time is less than two days on a single Tesla V100 GPU with our implementation in PyTorch [40]. We set $\lambda_{sty} = 1$, $\lambda_{ds} = 1$, and $\lambda_{cyc} = 1$ for CelebA-HQ and $\lambda_{sty} = 0.3$, $\lambda_{ds} = 1$, and $\lambda_{cyc} = 0.1$ for AFHQ. To stabilize the training, the weight λ_{ds} is linearly decayed to zero over the first 50K iterations. We adopt the non-saturating adversarial loss [10] with R_1 regularization [34] using $\gamma = 1$. We use the Adam [24] optimizer with $\beta_1 = 0$ and $\beta_2 = 0.99$. The learning rates for G , D , and E are set to 10^{-4} , while that of F is set to 10^{-6} . For evaluation, we employ exponential moving averages over parameters [21, 48] of all modules except D . We initialize the weights of all modules using He initialization [12] and set all biases to zero, except for the biases associated with the scaling vectors of AdaIN that are set to one.

¹<https://www.flickr.com>

²<https://www.pixabay.com>

C. Evaluation protocol

This section provides details for the evaluation metrics and evaluation protocols used in all experiments.

Frechét inception distance (FID) [14] measures the similarity between two sets of images. We use the feature vectors from the last average pooling layer of the ImageNet-pretrained Inception-V3 [43]. For each test image from a source domain, we translate it into a target domain using 10 latent vectors, which is randomly sampled from the standard Gaussian distribution. We then compute the FID between the translated images and training images in the target domain. We compute the FIDs for every pair of image domains (e.g., female \rightleftarrows male for CelebA-HQ) and report the average score. Note that, for reference-guided synthesis, each source image is transformed using 10 reference images randomly sampled from the test set of the target domain.

Learned perceptual image patch similarity (LPIPS) [51] measures the diversity of the generated images using the L1 distance between features extracted from the pretrained AlexNet [25]. For each test image from a source domain, we generate 10 output images of a target domain using 10 randomly sampled latent vectors. Then, we compute the average of the pairwise distances among all outputs generated from the same input (i.e., 45 pairs). Finally, we report the average of the LPIPS values over all test images. For reference-guided synthesis, each source image is transformed using 10 reference images to produce 10 outputs.

D. Additional results

We provide additional reference-guided image synthesis results on both CelebA-HQ and AFHQ (Figure 9 and 10). In CelebA-HQ, StarGAN v2 synthesizes the source identity in diverse appearances reflecting the reference styles such as hairstyle, makeup, and beard. In AFHQ, the result images follow the breed and hair of the reference images preserving the pose of the source images. Interpolation results between styles can be found in the accompanying videos.

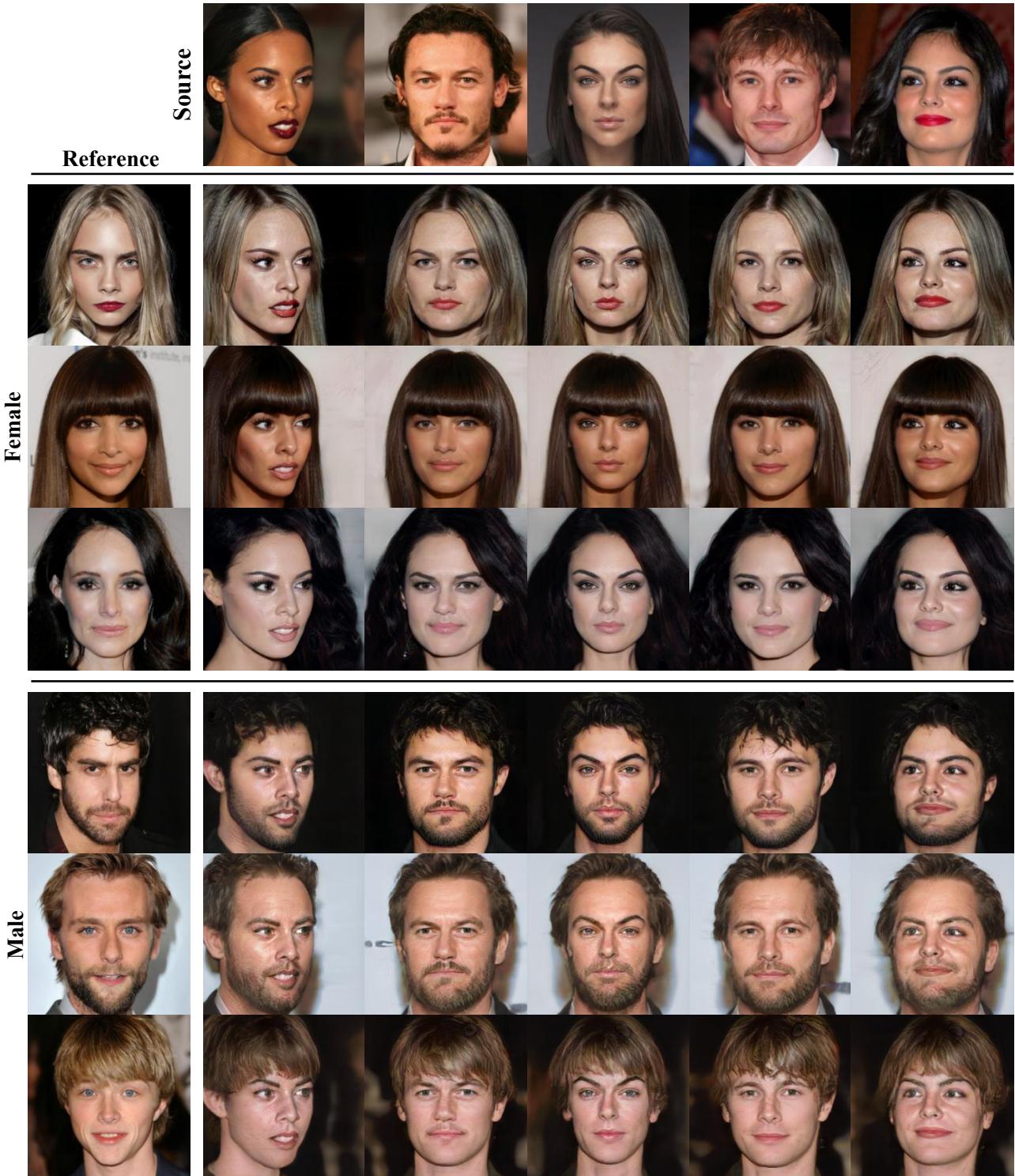


Figure 9. Reference-guided image synthesis results on CelebA-HQ. The source and reference images in the first row and the first column are real images, while the rest are images generated by our proposed model, StarGAN v2. Our model learns to transform a source image reflecting the style of a given reference image. High-level semantics such as hairstyle, makeup, and beard are followed from the reference images, while the pose and identity of the source images are preserved. Note that the images in each column share a single identity with different styles, and those in each row share a style with different identities.

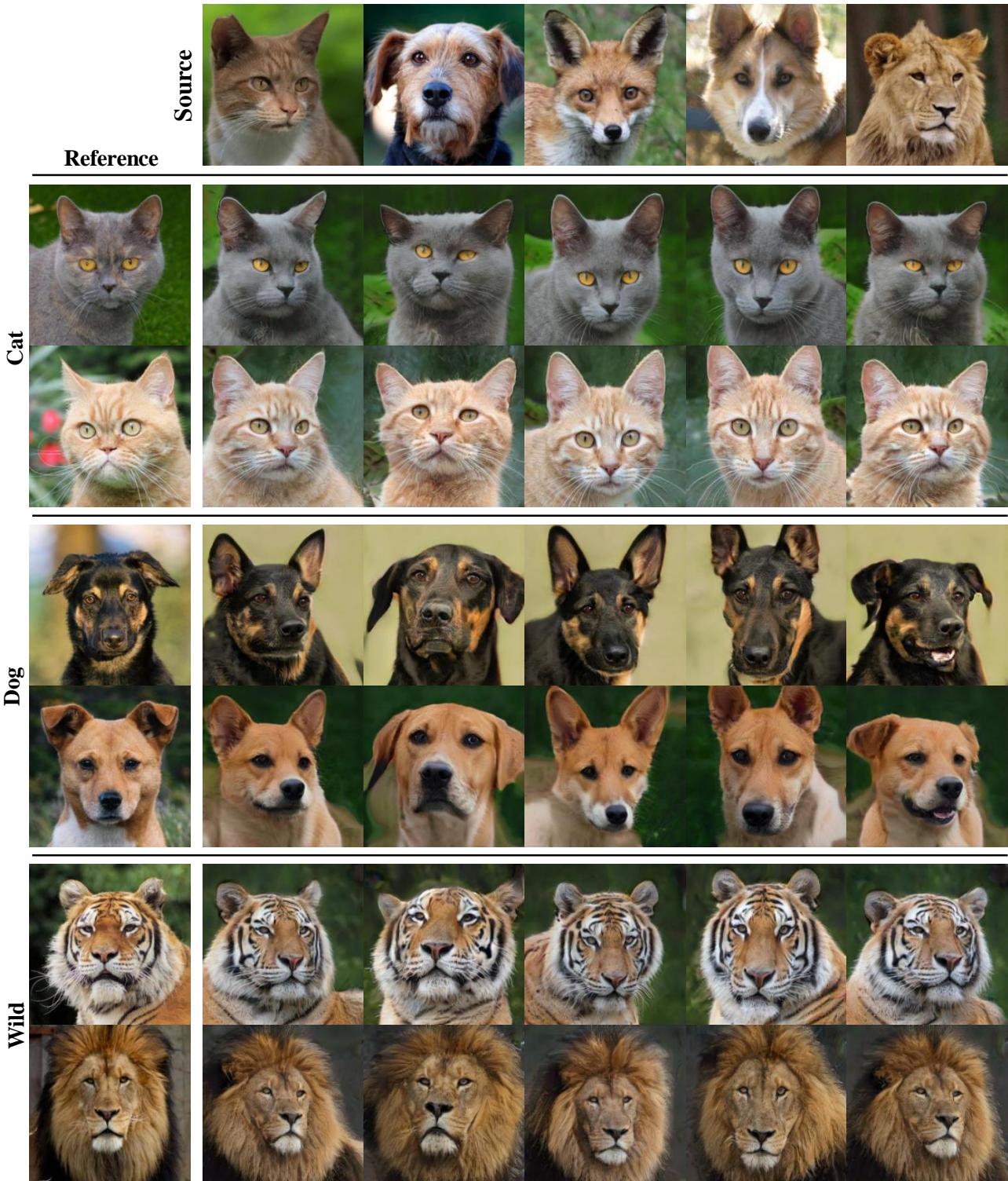


Figure 10. Reference-guided image synthesis results on AFHQ. The source and reference images in the first row and the first column are real images, while the rest are images generated by our proposed model, StarGAN v2. Our model learns to transform a source image reflecting the style of a given reference image. High-level semantics such as breed and hair texture are followed from the reference images, while the gaze and pose of the source images are preserved.

E. Network architecture

In this section, we provide architectural details of StarGAN v2, which consists of four modules described below.

Generator (Table 5). Our generator consists of four down-sampling blocks, four intermediate blocks, and four up-sampling blocks, all of which inherit pre-activation residual units [13]. We use the instance normalization (IN) [44] and the adaptive instance normalization (AdaIN) [15, 22] for down-sampling and up-sampling blocks, respectively. A style code is injected into all AdaIN layers, providing scaling and shifting vectors through learned affine transformations. We use the average pooling for down-sampling and the nearest-neighbor interpolation for up-sampling. We do not use the hyperbolic tangent as an output activation and let the model to learn the output color range.

Mapping network (Table 6). Our mapping network consists of an MLP with N output branches, where N indicates the number of domains. Seven fully connected layers are shared among all domains, followed by one specific fully connected layer for each domain. We set the dimensions of the latent code, the hidden layer, and the style code to 16, 512, and 64, respectively. We sample the latent code from the standard Gaussian distribution. We do not apply the pixel normalization [22] to the latent code, which has been observed not to increase model performance in our tasks. We also tried feature normalizations [3, 19], but this degraded performance.

Style encoder (Table 7). Our style encoder consists of a CNN with N output branches, where N is the number of domains. Six pre-activation residual blocks are shared among all domains, followed by one specific fully connected layer for each domain. We do not use the global average pooling [16] to extract fine style features of a given reference image. The output dimension “D” in the Table 7 is set to 64, which indicates the dimension of the style code.

Discriminator (Table 7). Our discriminator is a multi-task discriminator [34], which contains multiple linear output branches ³. The discriminator contains six pre-activation residual blocks with leaky ReLU [32]. We use N fully-connected layers for real/fake classification of each domain, where N indicates the number of domains. The output dimension “D” is set to 1 for real/fake classification. We do not use any feature normalization techniques [19, 44] nor PatchGAN [20] as they have been observed not to improve output quality. We have observed that in our settings, the multi-task discriminator provides better results than other types of conditional discriminators [35, 36, 38, 41].

³The original implementation of the multi-task discriminator can be found at https://github.com/LMescheder/GAN_stability.

LAYER	RESAMPLE	NORM	OUTPUT SHAPE
Image x	-	-	$256 \times 256 \times 3$
Conv1x1	-	-	$256 \times 256 \times 32$
ResBlk	AvgPool	IN	$128 \times 128 \times 64$
ResBlk	AvgPool	IN	$64 \times 64 \times 128$
ResBlk	AvgPool	IN	$32 \times 32 \times 256$
ResBlk	AvgPool	IN	$16 \times 16 \times 512$
ResBlk	-	IN	$16 \times 16 \times 512$
ResBlk	-	IN	$16 \times 16 \times 512$
ResBlk	-	AdaIN	$16 \times 16 \times 512$
ResBlk	-	AdaIN	$16 \times 16 \times 512$
ResBlk	Upsample	AdaIN	$32 \times 32 \times 256$
ResBlk	Upsample	AdaIN	$64 \times 64 \times 128$
ResBlk	Upsample	AdaIN	$128 \times 128 \times 64$
ResBlk	Upsample	AdaIN	$256 \times 256 \times 32$
Conv1x1	-	-	$256 \times 256 \times 3$

Table 5. Generator architecture.

LAYER	ACTVATION	NORM	OUTPUT SHAPE
Latent z	-	-	<u>16</u>
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear	ReLU	-	512
Linear * N	-	-	$64 * N$

Table 6. Mapping network architecture.

LAYER	RESAMPLE	NORM	OUTPUT SHAPE
Image x	-	-	$256 \times 256 \times 3$
Conv1x1	-	-	$256 \times 256 \times \text{ch}$
ResBlk	AvgPool	-	$128 \times 128 \times 2\text{ch}$
ResBlk	AvgPool	-	$64 \times 64 \times 4\text{ch}$
ResBlk	AvgPool	-	$32 \times 32 \times 8\text{ch}$
ResBlk	AvgPool	-	$16 \times 16 \times 16\text{ch}$
ResBlk	AvgPool	-	$8 \times 8 \times 32\text{ch}$
ResBlk	AvgPool	-	$4 \times 4 \times 32\text{ch}$
LReLU	-	-	$4 \times 4 \times 32\text{ch}$
Conv4x4	-	-	$1 \times 1 \times 32\text{ch}$
LReLU	-	-	$1 \times 1 \times 32\text{ch}$
Reshape	-	-	32ch
Linear * N	-	-	$D * N$

Table 7. Style encoder and discriminator architectures. “ch” represents the channel multiplier and is set to 16 and 32 for the style encoder and the discriminator, respectively.

References

- [1] A. Almahairi, S. Rajeshwar, A. Sordoni, P. Bachman, and A. Courville. Augmented cyclegan: Learning many-to-many mappings from unpaired data. In *ICML*, 2018. 2, 8
- [2] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool. Combogan: Unrestrained scalability for image domain translation. In *CVPRW*, 2018. 2
- [3] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. In *arXiv preprint*, 2016. 12
- [4] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019. 8
- [5] H. Chang, J. Lu, F. Yu, and A. Finkelstein. Pairedcyclegan: Asymmetric style transfer for applying and removing makeup. In *CVPR*, 2018. 8
- [6] W. Cho, S. Choi, D. K. Park, I. Shin, and J. Choo. Image-to-image translation via group-wise deep whitening-and-coloring transformation. In *CVPR*, 2019. 8
- [7] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 2, 3, 4
- [8] J. Donahue and K. Simonyan. Large scale adversarial representation learning. In *NeurIPS*, 2019. 8
- [9] V. Dumoulin, E. Perez, N. Schucher, F. Strub, H. d. Vries, A. Courville, and Y. Bengio. Feature-wise transformations. In *Distill*, 2018. 4
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. In *NeurIPS*, 2014. 8, 9
- [11] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein gans. In *NeurIPS*, 2017. 4
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015. 9
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 12
- [14] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 4, 9
- [15] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 2, 4, 12
- [16] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 2, 3, 4, 6, 7, 8, 12
- [17] L. Hui, X. Li, J. Chen, H. He, and J. Yang. Unsupervised multi-domain image translation with domain-specific encoders/decoders. In *ICPR*, 2018. 2
- [18] K. Hyunsu, J. Ho Young, P. Eunhyeok, and Y. Sungjoo. Tag2pix: Line art colorization using text tag with secat and changing loss. In *ICCV*, 2019. 8
- [19] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 12
- [20] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial nets. In *CVPR*, 2017. 1, 8, 12
- [21] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018. 4, 9
- [22] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2, 8, 12
- [23] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. 3
- [24] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 9
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 9
- [26] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 8
- [27] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. K. Singh, and M.-H. Yang. Diverse image-to-image translation via disentangled representations. In *ECCV*, 2018. 2, 3, 4, 6, 7, 8
- [28] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NeurIPS*, 2017. 8
- [29] M.-Y. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz. Few-shot unsupervised image-to-image translation. In *ICCV*, 2019. 2, 4, 8
- [30] M. Lucic, M. Tschannen, M. Ritter, X. Zhai, O. Bachem, and S. Gelly. High-fidelity image generation with fewer labels. In *ICML*, 2019. 8
- [31] L. Ma, X. Jia, S. Georgoulis, T. Tuytelaars, and L. Van Gool. Exemplar guided unsupervised image-to-image translation with semantic consistency. In *ICLR*, 2019. 8
- [32] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML*, 2013. 12
- [33] Q. Mao, H.-Y. Lee, H.-Y. Tseng, S. Ma, and M.-H. Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *CVPR*, 2019. 2, 3, 4, 6, 7, 8
- [34] L. Mescheder, S. Nowozin, and A. Geiger. Which training methods for gans do actually converge? In *ICML*, 2018. 2, 4, 8, 9, 12
- [35] M. Mirza and S. Osindero. Conditional generative adversarial nets. In *arXiv preprint*, 2014. 4, 12
- [36] T. Miyato and M. Koyama. cGANs with projection discriminator. In *ICLR*, 2018. 12
- [37] S. Na, S. Yoo, and J. Choo. Miso: Mutual information loss with stochastic style representations for multimodal image-to-image translation. In *arXiv preprint*, 2019. 2
- [38] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. In *ICML*, 2017. 4, 12
- [39] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019. 8

- [40] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NeurIPS*, 2017. 9
- [41] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 12
- [42] A. Romero, P. Arbeláez, L. Van Gool, and R. Timofte. Smit: Stochastic multi-label image-to-image translation. In *ICCVW*, 2019. 4
- [43] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 9
- [44] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. In *arXiv preprint*, 2016. 12
- [45] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *ECCV*, 2018. 8
- [46] Y. Wang, A. Gonzalez-Garcia, J. van de Weijer, and L. Herranz. Sdit: Scalable and diverse cross-domain image translation. In *ACM-MM*, 2019. 4
- [47] D. Yang, S. Hong, Y. Jang, T. Zhao, and H. Lee. Diversity-sensitive conditional generative adversarial networks. In *ICLR*, 2019. 3
- [48] Y. Yazici, C.-S. Foo, S. Winkler, K.-H. Yap, G. Piliouras, and V. Chandrasekhar. The unusual effectiveness of averaging in gan training. In *ICLR*, 2019. 9
- [49] S. Yoo, H. Bahng, S. Chung, J. Lee, J. Chang, and J. Choo. Coloring with limited data: Few-shot colorization via memory augmented networks. In *CVPR*, 2019. 8
- [50] X. Yu, Y. Chen, T. Li, S. Liu, and G. Li. Multi-mapping image-to-image translation via learning disentanglement. In *NeurIPS*, 2019. 8
- [51] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 4, 9
- [52] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 3, 8
- [53] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman. Toward multimodal image-to-image translation. In *NeurIPS*, 2017. 2, 8