

18

19

20

21

22

23

24

# Article Masked Style Transfer for Source-coherent Image-to-Image Translation

Filippo Botti \*🝺, Tomaso Fontanini 🖻, Massimo Bertozzi 🖻, Andrea Prati 🖻

Department of Engineering and Architecture, University of Parma, 43124 Parma, Italy \* Correspondence: filippo.botti@unipr.it

Abstract: The goal of image-to-image translation (I2I) is to translate images from one domain to 1 another while maintaining the content representations. A popular method for I2I translation involves the use of a reference image to guide the transformation process. However, most of the architectures 3 fail to maintain the input main characteristics and produce images too close to the reference during 4 style transfer. In order to avoid this problem, we propose a novel architecture which is able to perform source-coherent translation between multiple domains. Our goal is to preserve input details during 6 I2I translation by weighting the style code obtained from the reference images before applying it to the source image. Therefore, we choose to mask the reference images in an unsupervised way, before extracting the style from them. By doing so, the input characteristics while performing style 9 transfer are better maintained. As a result, we also increase the diversity in images generated by 10 extracting style from the same reference. Additionally, the adaptive normalization layers, commonly 11 used to inject styles into the model, are substituted with an attention mechanism for the purpose of 12 increasing the quality in generated images. Several experiments are performed on CelebA-HQ and 13 AFHQ datasets in order to prove the efficacy of the proposed system. Quantitative results, measured 14 with LPIPS and FID metrics, demonstrate the superiority the proposed architecture compared to the 15 state of art. 16

Keywords: deep learning; style transfer; image-to-image translation; generative adversarial networks

## 1. Introduction

Image-to-image translation (I2I) aims to generate an output image with a different style while preserving the content information of the input [1]. More specifically, the goal of I2I is to convert an image  $x_A$ , belonging to a *source* domain A, into an image  $y_B$ , belonging to a *target* domain B, by preserving its intrinsic contents belonging to the source domain and modifying its extrinsic contents by making them as similar as possible to those characterizing the target domain.

A lot of frameworks that use generative models to perform I2I translation are emerging 25 in a variety of areas: from face editing [2], to style transfer [3] and automotive field [4]. 26 Focusing on style transfer, StarGANv2 [5] introduced an innovative approach. Specifically, 27 StarGANv2 incorporates a Style Encoder, designed to extract the style characteristics of 28 an image referred to as the *reference image*. Subsequently, the extracted style is applied to 29 the input image using a single Generator that is able to perform image translation across 30 multiple domains. StarGANv2 also has a Mapping Network in charge of generating styles 31 for the Generator from random noise. In their work, authors of StarGANv2 introduce 32 diversity as the characteristic of each image within a domain to be different despite coming 33 from the same domain. By this definition, the authors show how the output changes as 34 the reference changes, even if picked from the same domain. However, this architecture 35 design tends to apply global changes to the entire input image without preserving its 36 intrinsic content representation. This can be described as a heavy form of reference-based 37 style transfer, which can be seen in the output images generated by extracting the style from 38

Citation: Botti, F.; Fontanini, T.; Bertozzi, M.; Prati, A. Masked Style Transfer for Source-coherent Image-to-Image Translation. *Appl. Sci.* 2024, 1, 0. https://doi.org/

Received: Revised: Accepted: Published:

**Copyright:** © 2025 by the author. Submitted to *Appl. Sci.* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). the same reference. In such cases, the output images tend to closely mirror the reference image and the generated images collapse to the reference image as it can be seen in Fig. 1.

To address this limitation, we present a novel architecture which is able to perform 41 source-coherent I2I translation across multiple domains. Our solution involves adding a 42 segmentation layer before the Style Encoder: this layer computes segmentation masks 43 which will be used to separate the subjects within the reference image and to select only the 44 desired part of the image. In this way, we can remove all of the unnecessary content in the 45 reference, like background or out-of-domain parts. Ultimately, the Style Encoder extracts 46 the style only from the relevant part of the reference image and the Generator produces 47 images with the style of the reference image without collapsing to it. 48

Our network architecture takes inspiration from StarGANv2 [5], though with some fundamental changes; in particular, we change the style application by using Cross Attention layers [6] and not Adaptive Instance Normalization (AdaIN) [7]; then, we adapted the Style Encoder by feeding it with both image and its correspondent mask.

Moreover, a crucial aspect lies in the utilization of an unsupervised architecture for extracting masks from reference images. Specifically, we choose to use STEGO [8] model, an architecture which can perform unsupervised semantic segmentation, in order to produce masks. With STEGO we produce binary images, which will be used in order to separate information on where to extract the style code inside the Style Encoder.

To summarize, the main contributions of the proposed work are the followings:

- **Innovative architecture for style transfer.** Introduction of a novel architecture which is able to perform source-coherent I2I translation between multiple domains by preserving input details and increasing diversity during generation.
- Semantic style separation. The model utilizes unsupervised segmentation architecture to produce masks in order to localize the style only on specific subjects of the images and removing useless areas like background or out-of-domain details. By this weighing of the reference images, the model is able to focus only on the relevant parts and better understand the characteristics of the style images, resulting in more accurate style code compared to the ones generated by state-of-art architectures.
- **Transferring style using Cross Attention.** The proposed architecture also shows how attention mechanisms, more in details Cross-Attention layers, are able to improve the quality of style transfer, with respect to commonly used Adaptive Instance Normalization layers.

#### 2. Related work

**Image-to-image translation.** Image-to-image translation was first introduced by [9] as 73 the task of translating one possible representation of a scene into another, given sufficient 74 training data. Pix2Pix [10] was the first attempt to use GANs, in particular Conditional 75 GANs (CGANs), in order to translate an image from source domain to target domain and 76 viceversa with paired datasets. Later, CycleGAN [3] improved Pix2Pix performance by 77 removing the requirement of paired datasets and suggested a method for I2I translation 78 on unpaired datasets by employing a cycle consistency loss that guarantees that an image 79 should accurately replicate the source image when it is translated to the target domain 80 and then reversed. MUNIT [11] was one of the first attempt to enhance diversity of the 81 generated images by feeding the generator with a style code that is randomly sampled 82 from Gaussian noise. Later, MSGAN [12], tried to improve diversity of generated images 83 by maximizing the ratio of the distance of two images in the image space with respect to 84 the distance of their correspondent latent code in the latent space. StarGAN [13] reached 85 better performance both on diversity and quality terms by using only a single generator to train between multiple domains. StarGANv2 [5] later improved StarGAN architecture by 87 introducing a Style Encoder which is in charge to learn style from image and then use this 88 style code in order to condition the output. Nevertheless, all of the cited architectures tend 89 to share the same limitation of lack of diversity when using the same reference. 90

39

40

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

97

123

Diffusion Probabilistic Models (DPM) [14] has recently showed impressive results in the generative field. Despite this, DPMs are still not at the same level of GANs for I2I translation problems. Architectures like ControlNet [15], BBDM [16] or Palette [17] show good results for primitive form of I2I, but they lack of capacity to perform I2I from multiple domains. For this reason, in this paper, we chose to adapt StarGANv2 architecture to perform our task.

Style transfer. Style transfer is a way to perform I2I translation by generating a sample with 98 the same content of the input image but with another style. In this way, we can translate qq images between multiple domains and preserving the input intrinsic characteristics. One of 100 the first application used conditional GAN [18] in order to perform style transfer, but was 101 based on a slow optimization process that iteratively updates the image to minimize content 102 and style losses. Later, Adaptive Instance Normalization (AdaIN) [7] became the state of 103 art for style transfer application. AdaIN enables fast arbitrary style transfer in real-time 104 without being limited to a specific set of styles as in previous works. Recently, Transformers 105 [19] exhibited impressive results in NLP and a lot of Transformer-based architecture have 106 been used across a multitude of vision-related tasks. In particular, StyTr<sup>2</sup> [20] and Latent 107 Diffusion [6] model highlighted the power of Transformers and Cross-Attention layers 108 when used to transfer style from multimodal reference like text, class labels or images. 109 For this reason, we selected Cross Attention layers in order to apply domain style to 110 the generated input. Additionally, recent approaches leveraged the capability of Latent 111 Diffusion to perform style transfer between pictures and paintings [21,22]. 112

One of the main challenges during style transfer is to identify only the regions where 113 to extrapolate the style and remove unnecessary regions like background or other parts 114 of the image. [23] introduced an attention layer in order to select the area on where to 115 apply style during I2I translation. [24] proposed a cycle consistent attention loss in order to 116 train the model to apply changes on the same area during translation and reconstruction 117 by using residual block activation map. Recently, SEAN [25] demonstrated that using a 118 mask that represents only the relevant area of the image is possible to perform an average 119 pooling operation on the extracted features inside the style encoder and to produce more 120 accurate style codes. Following a similar idea, we propose to modify StarGANv2 style 121 encoder by introducing the mask multiplication and pooling. 122

Unsupervised semantic segmentation. Semantic segmentation aims to discover and local-124 ize semantically meaningful categories present in an image. Tipically, Mask R-CNN [26] 125 or YOLO [27] are used in order to produce segmentation from an image, but they require 126 labelled datasets and this is not always feasible and, in any case, not scalable. Recently, sev-127 eral works introduced semantic segmentation systems that could learn from weaker forms 128 of labels, such as classes, tags, bounding boxes, scribbles, or point annotations. IIC system 129 [28] focuses on maximizing the mutual information of patch-level cluster assignments 130 between an image and its augmentations. It operates as an implicit clustering method, with 131 the network directly predicting the (soft) clustering assignment for each pixel-level feature 132 vector [8,28,29]. PiCIE [29] enhances the semantic segmentation outcomes achieved by IIC 133 by leveraging invariance to photometric effects and equivariance to geometric transforma-134 tions as an inductive bias. In PiCIE, the network aims to minimize the distance between 135 features subjected to different transformations. The distance metric is determined through 136 an in-the-loop k-means clustering process [8,29]. Conversely, STEGO [8] gains impressive 137 results in semantic segmentation without any kind of labelled dataset. STEGO shows that 138 unsupervised deep network features have correlation patterns that are largely consistent 139 with true semantic labels and uses these patterns to categorize every pixel of the image. 140 Based on its valuable characteristics, STEGO is a perfect candidate for our purposes and it 141 represents the state of art in unsupervised semantic segmentation. 142





**Figure 1.** Results generated using StarGANv2. It can be seen how the generated images lose input intrinsic characteristics like hair length or coat color and end up looking too similar to the reference, resulting in a lack of diversity.

#### 3. Proposed system

In the next sections, the proposed model, which is able to perform source-coherent translation by preventing results to collapse to the references, will be described.

#### 3.1. Network Architecture

As stated above, the network architecture of our system follows that of StarGANv2 147 and is composed by a Generator G, a Discriminator D, a Style Encoder E and a Mapping 148 Network *M* (see Fig. 2). Style code  $s_{trg}$  can be generated both from images by using the 149 Style Encoder or from random noise z by using the Mapping Network. During generation, 150 the Generator *G* takes both source image  $x_{src}$  and style code  $s_{trg}$  and generates the output 151  $x_{trg} = G(x_{src}, s_{trg})$ . In finer details, G is designed as an encoder-decoder architecture 152 featuring four downsampling residual blocks and four upsampling residual blocks, but 153 uses Cross-attention layers (instead of Adaptive Instance Normalization layers) to apply 154 style. The discriminator D serves the role of evaluating which domain the generated 155 samples belong to. D follows StarGANv2 implementation, and is a multitask discriminator, 156 which consists of multiple output branches. Finally, the style encoder *E* is a CNN with two 157 residual blocks as features extractor and an average pooling layer on the area covered by 158 the mask, while *M* is an MLP in charge of generating style codes from noise.

#### 3.1.1. Generating Style from Reference

As previously stated, the Style Encoder is heavily modified w.r.t. the one introduced 161 by StarGANv2. In fact, we took inspiration from SEAN Style Encoder [25] and adapted 162 it to our goal. First, only two, instead of six, downsampling layers are used, following 163 the implementation reported in [30], because we do not need shape information inside 164 our style code. Then, we multiply the extracted features with the precomputed mask, 165 computed from segmentation map obtained with STEGO (see section 3.3). This allows us 166 to delete information irrelevant for the style computation, such as the background. Finally, 167 an average pooling followed by a fully connected layer for each domain is applied (see Fig. 168 3(a)). 169

#### 3.2. Transferring Style with Cross Attention

In order to apply style to the image, we use Cross-Attention layers instead of AdaIN during the decoding phase of the Generator. As shown in Fig. 3(b), we first extract features from source image  $x_{src}$  and style code  $s_{trg}$  extracted from  $x_{ref}$ . Subsequently, the features 173

145 146

143

144

159 160



**Figure 2.** Overview of proposed system. The generator takes the input image  $x_{src}$  and applies the style code, computed using reference  $(x_{ref})$  and its correspondent mask  $(m_{ref})$  or using random noise (z), with Cross Attention layers. Finally, the result  $x_{trg}$  passes through the Discriminator.

extracted from the source image are normalized with a Layer Normalization and then, for 174 every Cross-Attention, the style code is injected as follows: 175

$$Att(x,s) = Softmax\left(\frac{Q\cdot K^{T}}{\sqrt{d}}\right) \cdot V$$
<sup>176</sup>

where *Q* is the projection of the features extracted from  $x_{src}$ , *K*, *V* are projections of the style 177 code  $s_{trg}$  and d is the dimension of a single attention head. Finally, the resulting tensors are 178 normalized with Layer Normalization and linearly transformed with a Feed Forward layer. 179

#### 3.3. Extracting Mask with STEGO

For the purpose of maintaining the network fully unsupervised, the unsupervised 181 semantic segmentation architecture STEGO is employed for extracting masks from the 182 reference image before extracting style from it. As described in [8], features f are first 183 extracted from the reference image using DiNo [31] feature extractor, then STEGO segmen-184 tation head is devoted to extract a non-linear projection and to learn patterns inside the 185 image. Finally, results are clustered and refined with a Conditional Random Field (CRF) 186 layer [32]. 187

Only the semantic cluster of the style that needs to be transferred is selected and the 188 others are set to zero in the segmentation mask. More in detail, the selected semantic 189 cluster/class is the one corresponding to the main subject of the image (e.g., animals in 190 AFHQ dataset [5], person in CelebA-HQ dataset [33]). 191

### 3.4. Training and Losses

In order to train the proposed model, we choose to maintain the training phase of StarGANv2 without any change. Therefore, the total loss is composed of four losses:

Adversarial loss used to learn the generation of realistic results:

$$\mathcal{L}_{adv} = \mathbb{E}_{src}[\log D_{src}(x_{src})] + \mathbb{E}_{trg}[\log(1 - D_{trg}(G(x_{src}, s_{trg})))]$$

where  $x_{src} \in \mathcal{X}$  is the input image,  $G(\cdot)$  is the generator that takes  $x_{src}$  and  $s_{trg}$ , which 197 is the style code extracted from the reference image  $x_{ref}$ .

180

192

193

194

195

196



Figure 3. Overview of Style Encoder architecture. (a) Comparison between StarGANv2 Style Encoder (left) and the proposed Style Encoder (right). Differently from StarGANv2, we utilize Mask for computing the style code. (b) Overview of Cross Attention layers utilized for style transfer. We compute Cross Attention using extracted features of the input image as Query vector and the style code as Key and Value vector.

Style Reconstruction loss introduced in order to avoid the generator G from ignoring 199 the style  $s_{trg}$  during the generation phase: 200

$$\mathcal{L}_{sty} = \mathbb{E}_{src,trg} \left[ \left\| s_{trg} - E(x_{trg}) \right\| \right]$$
<sup>201</sup>

where  $E(x_{trg})$  is the style code extracted from the generated image.

Style Diversification loss used to differentiate the style generated from two different images:

$$\mathcal{L}_{div} = \mathbb{E}_{src,trg_1,trg_2} \left[ \left\| x_{trg_1} - x_{trg_2} \right\| \right]$$
<sup>20</sup>

*Cycle Consistency loss* to maintain the domain-invariant characteristics of the generated image like pose and shape:

$$\mathcal{L}_{cyc} = \mathbb{E}_{src,trg} \left[ \left\| x_{src} - G(x_{trg}, \tilde{s}_{src}) \right\| \right]$$
<sup>208</sup>

where  $\tilde{s}_{src}$  is the estimated style code extracted from the input. The final loss is therefore as follows:

$$\min_{G,M,E} \max_{D} \mathcal{L}_{adv} + \lambda_{sty} \mathcal{L}_{sty} - \lambda_{div} \mathcal{L}_{div} + \lambda_{cyc} \mathcal{L}_{cyc}$$
<sup>211</sup>

It is worth noticing that, during training, a reference image and random noise are used 212 alternatively for generating the style code through the mapping network.

#### 4. Experimental Results

This section will report details on the experimental results, both qualitative and quantitative. 216

### 4.1. Selected Baseline

Since our work is an extension of StarGANv2, we decided to not compare the results 218 that we produce with other architectures, following [30] comparison. In fact, our work can 219 be seen as an improved version of StarGANv2 with the objective to show how to leverage 220 the mistakes made by StarGANv2 and how to improve them by adding our masked style 221 encoder. As discussed, StarGANv2 style transfer is limited in terms of diversity and tends 222 to generate images with the same style applied, on the contrary we show a proper way 223

214

213

202

203

204

206

207

209

210

215

to perform style transfer without losing diversity. Moreover, is quite difficult to compare our architecture with others since, when considering diffusion models for example, the style transfer is a totally different task and cannot be compared to our architecture. All the hyperparameters and training strategies are the ones proposed in the original paper of StarGANv2. 228

## 4.2. Datasets

We tested our model on two datasets: CelebA-HQ, composed by 30k images [33] and 230 AFHQ, composed by 16k images [5]. CelebA-HQ is organized in two domains (male and 231 female) and AFHQ in three (cat, dog and wildlife animals). Binary masks are extracted 232 using STEGO pretrained on COCOstuff [34] and selecting "person" and "animal" attributes 233 in order to identify the subject of the image for the two datasets. No other information is 234 employed during training or inference. We resized all images to  $256 \times 256$  and all masks to  $244 \times 64$  during training. 236

## 4.3. Implementation Details

During all the experiments, we trained the network for 100k iterations and we use 238 Adam [35] as optimizer. Learning rates of  $10^{-4}$  for *G*, *D*, *E* and  $10^{-6}$  for *M* are used. 239 Training took about 1 day on a single NVIDIA A100 GPU which is the same as the one 240 for StarGANv2 proving that our approach does not add complexity in the training. For 241 CelebA-HQ training, we weighed every loss equally; on the contrary, for AFHQ we set 242  $\lambda_{div}$  to 2 while  $\lambda_{cyc}$  and  $\lambda_{sty}$  to 1, following [5] implementation, in order to make an equal 243 comparison with StarGANv2 and to show that the results obtained are better because of 244 the architecture and not because of these hyperparameters. 245

#### 4.4. Evaluation metrics

In order to evaluate our model we used Frechét Inception Distance (FID) [36] for image quality and Learned Perceptual Image Patch Similarity (LPIPS) [37] to measure diversity in the generated results. More in detail, FID metric measures the distance between two distribution and in our case is used in order to measure the distance between generated images, e.g. cat generated images, and test set that contains real images, e.g. real cat images. So, intuitively, a low value of FID means that two distributions are similar. Indeed, given two Gaussian distribution (m, C) and  $(m_w, C_w)$  the FID is computed as follow: 249

$$d^{2}((m,C),(m_{w},C_{w})) = ||m-m_{w}||_{2}^{2} + \operatorname{Tr}(C+C_{w}-2(CC_{w})^{\frac{1}{2}})$$
<sup>254</sup>

The Learned Perceptual Image Patch Similarity (LPIPS) calculates perceptual similarity <sup>255</sup> between two images. LPIPS essentially computes the similarity between the activations <sup>256</sup> of two image patches for some pre-defined and pre-trained network. This measure has <sup>257</sup> been shown to match human perception well. A low LPIPS score means that image patches <sup>258</sup> are perceptual similar. Indeed, given two patches x and  $x_0$  their distance is computed as <sup>259</sup> follow: <sup>260</sup>

$$d(x, x_0) = \sum_{l} \frac{1}{H_l W_l} \sum_{h, w} ||w_l \odot (\hat{y}_{hw}^l - \hat{y}_{0hw}^l)||_2^2$$
<sup>263</sup>

where  $\hat{y}_{hw}^l$  and  $\hat{y}_{0hw}^l$  are the stacked features extracted from the patches. These features are normalized, and the distance between them is modulated by a learned weight vector  $w_l$  which adjusts the contribution of different feature channels [37].

Since the main contribution of our model is to perform source-coherent translation, which aims at improving diversity for image generated with the same reference, the evaluation was designed as follows: 267

- firstly, we randomly select one image for every domain as reference;
- secondly, given a set of source images, we generate samples with those reference images;
- thirdly, we compute FID and LPIPS (with consecutive pair of images);

229

236 237

246

268

finally, we repeat this evaluation phase for 10 times in order to remove randomness in 272 the results. 273

It is worth emphasizing that we decided to not use StarGANv2 FID algorithm, which 274 calculates FID by using ten references for each domain, because we want to improve 275 diversity in the results generated from a single reference. Therefore, FID is computed with 276 only one reference per domain in order to evaluate the quality of images generated with 277 our method. 278



Figure 4. Main results obtained with our architecture. We can see how the results are indistinguible from real images and how they maintain the input intrinsic characteristics like expression, age or fur color.

#### 4.5. Discussion

Depending on the dataset different styles are transferred. For CelebA-HQ the male2female 280 and *female2male* where chosen. While for the AFHQ dataset we transferred *cat-dog2wildlife* 281 and *wildlife2cat-dog*. As it can be seen from Fig. 4, the proposed architecture can perform 282 I2I translation between multiple domains like StarGANv2, but gaining the capability to 283 preserve input intrinsic characteristics during translation. Looking at the CelebA-HQ 284 results, the proposed architecture maintains the input facial attributes, but applies changes 285 to gender and hair color, which are taken from the reference images. In AFHQ results, our 286 method maintains the same expression and preserves better fur color during translation, 287 but changes the class of the animal. We introduced our work with the claim of increasing 288 diversity in results generated using the same reference. This is shown clearly in Figs. 5 289 and 6, where the results are compared with the ones obtained using StarGANv2. From 290 these examples, it is evident that StarGANv2 tends to collapse to the reference image and 291 looses the majority of intrinsic attributes from input, except for the pose in AFHQ and the 292 expression in CelebA-HQ. In contrast, our results maintain much more original details, 293 such as fur color in AFHQ or age and hair style in CelebA-HQ. This leeds to more diversity 294 and variety in generated images and in less reference-based generation. 295

More in details, StarGANv2 seems to capture attributes from reference, like hair 296 color and length in Fig. 5, and apply them to the input in a rigid scheme that does not 297 take care of the input attributes. On the contrary, the proposed method, in addition to 298 understanding which are the main details from the reference, also considers the details 299 from the input image before applying it. This leads, for instance, to a higher variety of 300 hair lengths and not in a prefixed hair length, like StarGANv2. This is also visible in Fig. 6, 301 where StarGANv2 produces the same animal with different poses. Our method, on the 302 contrary, better understands the input characteristics and generates different breeds of 303 dog/cats/wildlife animals based on the input breed. 304

All the previous considerations are valid also when the style code  $s_{trg}$  is sampled from random noise by using the Mapping Network *M*. This is presented in Fig. 7, where our architecture produces various and more source-coherent results than the ones generated by StarGANv2.

In order to support our claim we also show how our method produces similar results when similar reference images are employed, as shown in Fig. 8. This can be seen as a positive effect of our source-coherent method which does not ignore input attributes.

## 4.5.1. Quantitative results

Architecture	AFHQ		CelebA-HQ	
	$FID\downarrow$	LPIPS $\uparrow$	$FID\downarrow$	LPIPS $\uparrow$
StarGANv2 [2]	104.86	0.457	81.175	0.365
Ours (AdaIN)	76.15	0.523	57.67	0.420
Ours	67.72	0.517	54.12	0.425

**Table 1.** Quantitative comparison between StarGANv2 and our architecture. For our architecture we also include the model with AdaIN instead of Cross Attention Layers.

The above considerations are reflected in the quantitative results reported in Table 1. <sup>313</sup> The proposed architecture significantly improves LPIPS results for both datasets. Furthermore, FID results highlight how our architecture produces much higher quality images. <sup>315</sup>



**Figure 5.** Comparison between StarGANv2 (first rows) and our architecture (second rows) on CelebA-HQ dataset.

Figure 6. Comparison between StarGANv2 (first rows) and our architecture (second rows) on AFHQ dataset.

Architecture	CelebA-HQ FID↓
StarGANv2 [2] Ours (AdaIN)	<b>29.88</b> 32.94
Ours	30.99

**Table 2.** Quantitative comparison between StarGANv2 and our architecture using StarGANv2 FID algorithm.

As shown in Table 2, we also compute FID using StarGANv2 algorithm on CelebA-HQ and we obtained opposite results. This is due to how FID works (also explained in Section 4.4): given the alignment of StarGANv2 generated images with the reference image shown in qualitative results, FID computed as in StarGANv2 original paper is natively lower. Indeed, FID tends to measure the difference between two distribution and since the images generated by StarGANv2 have less diversity than the ones generated by our architecture, the FID score in this case is better for StarGANv2 results. Nevertheless, in section 4.4 we 322





**Figure 7.** Comparison between StarGANv2 (first rows) and our architecture (second rows) on AFHQ dataset in results generated using random noise as a reference.



**Figure 8.** Similar input generates similar output due to the fact that we preserve input characteristics during translation.

justified how is not fair to compute FID in this way in order to consider the diversity in generated images. However, this comes at the cost of more limitations for StarGANv2 w.r.t. to our architecture, such as loosing input characteristics, lack of diversity in generated results and results collapsed to the reference images.

#### 4.6. Ablation

Finally, we perform ablation studies to find the optimal configuration for our archi-328 tecture. First, we try to transfer the style using AdaIN and not Cross Attention layers. 329 As shown in Table 1 and Fig. 9, using AdaIN leads to slightly better results compared 330 to StarGANv2 in terms of diversity, but the network still does not maintain the input 331 characteristics like our final configuration. Additionally, we also tested employing 2 and 332 3 downsampling layers inside our Style Encoder, as it can be seen from second and third 333 rows in Fig. 9. Indeed, the configuration with 3 downsampling layers tends to collapse 334 more to the reference than the ones with 2 downsampling layer, as it can be seen from the 335 fur style. Furthermore, by comparing these results with the ones produced by the final 336 architecture, it is evident how Cross Attention layers improve the quality in generated 337 results and better maintain input characteristics. Finally, results generated without masks 338 are reported in the fourth row, proving that masks are necessary to identify major input 339 information, like fur color or ear pose. 340



**Figure 9.** Difference between results generated with different architectures: first row with StarGANv2, second with our Style Encoder and AdaIN for style transfer, third with the same architecture as before but with 3 downsampling layers inside the Style Encoder, fourth with Cross Attention layers for style transfer but without masks and with 3 downsampling layers in the Style Encoder and finally with the proposed architecture.

#### 5. Conclusions

For the task of I2I translation, StarGANv2 has shown limitations in preserving input 342 details during translation. Additionally, StarGANv2 is not able to generate diverse samples 343 when using the same reference image. For these reasons, this paper proposes a novel 344 architecture for source-coherent image-to-image translation, which preserves input char-345 acteristics and increases diversity in the generated results. More in detail, style extracted 346 from the reference images are masked in order for the model to focus only over the relevant 347 information and are injected in the model using cross attention layers. By doing so, we 348 managed to improve both quantitative and qualitative results. 349

Future works could focus on improving results generated using two different references and only one source, since our method, by preserving input intrinsic characteristics, tends to produce similar results when the same source image is utilized.

Author Contributions: Conceptualization, F.B. and T.F.; methodology, F.B. and T.F.; software, F.B.;validation, T.F.; formal analysis, F.B. and T.F.; investigation, F.B.; resources, F.B.; data curation, F.B.;writing—original draft preparation, F.B.; writing—review and editing, F.B. and T.F. and M.B. andA.P.; visualization, F.B.; supervision, T.F. and M.B. and A.P.; project administration, F.B. and T.F. andM.B. and A.P. All authors have read and agreed to the published version of the manuscript.

Informed Consent Statement: Not Applicable.

 Data Availability Statement: Both CelebA-HQ and AFHQ datasets can be downloaded from: https:
 359

 //github.com/clovaai/stargan-v2 (accessed on 23 July 2024).
 360

Conflicts of Interest: The authors declare no conflict of interest.

341

358

### References

- Pang, Y.; Lin, J.; Qin, T.; Chen, Z. Image-to-Image Translation: Methods and Applications. *IEEE Transactions on Multimedia* 2022, 24, 3859–3881. https://doi.org/10.1109/TMM.2021.3109419.
- 2. Li, Y.A.; Zare, A.; Mesgarani, N. Starganv2-vc: A diverse, unsupervised, non-parallel framework for natural-sounding voice conversion. *arXiv preprint arXiv:*2107.10394 **2021**.
- Zhu, J.; Park, T.; Isola, P.; Efros, A.A. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. CoRR 2017, abs/1703.10593, [1703.10593].
- 4. Zheng, Z.; Wu, Y.; Han, X.; Shi, J. ForkGAN: Seeing into the Rainy Night. In Proceedings of the The IEEE European Conference on Computer Vision (ECCV), August 2020.
- 5. Choi, Y.; Uh, Y.; Yoo, J.; Ha, J.W. Stargan v2: Diverse image synthesis for multiple domains. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 8188–8197.
- 6. Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; Ommer, B. High-resolution image synthesis with latent diffusion models. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 10684–10695.
- 7. Huang, X.; Belongie, S. Arbitrary style transfer in real-time with adaptive instance normalization. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2017, pp. 1501–1510.
- 8. Hamilton, M.; Zhang, Z.; Hariharan, B.; Snavely, N.; Freeman, W.T. Unsupervised Semantic Segmentation by Distilling Feature Correspondences, 2022, [arXiv:cs.CV/2203.08414].
- 9. Isola, P.; Zhu, J.Y.; Zhou, T.; Efros, A.A. Image-to-image translation with conditional adversarial networks. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1125–1134.
- 10. Mirza, M.; Osindero, S. Conditional Generative Adversarial Nets. CoRR 2014, abs/1411.1784, [1411.1784].
- 11. Huang, X.; Liu, M.Y.; Belongie, S.; Kautz, J. Multimodal unsupervised image-to-image translation. In Proceedings of the Proceedings of the European conference on computer vision (ECCV), 2018, pp. 172–189.
- 12. Mao, Q.; Lee, H.Y.; Tseng, H.Y.; Ma, S.; Yang, M.H. Mode seeking generative adversarial networks for diverse image synthesis. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 1429–1437.
- 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 Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 8789–8797.
- 14. Ho, J.; Jain, A.; Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems* **2020**, 33, 6840–6851.
- 15. Zhang, L.; Rao, A.; Agrawala, M. Adding conditional control to text-to-image diffusion models. In Proceedings of the Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 3836–3847.
- 16. Li, B.; Xue, K.; Liu, B.; Lai, Y.K. Bbdm: Image-to-image translation with brownian bridge diffusion models. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition, 2023, pp. 1952–1961.
- 17. Saharia, C.; Chan, W.; Chang, H.; Lee, C.; Ho, J.; Salimans, T.; Fleet, D.; Norouzi, M. Palette: Image-to-image diffusion models. In Proceedings of the ACM SIGGRAPH 2022 conference proceedings, 2022, pp. 1–10.
- 18. Gatys, L.A.; Ecker, A.S.; Bethge, M. Image style transfer using convolutional neural networks. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2414–2423.
- 19. Vaswani, A. Attention is all you need. Advances in Neural Information Processing Systems 2017.
- 20. Deng, Y.; Tang, F.; Dong, W.; Ma, C.; Pan, X.; Wang, L.; Xu, C. Stytr2: Image style transfer with transformers. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 11326–11336.
- Chung, J.; Hyun, S.; Heo, J.P. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 8795–8805.
- 22. Deng, Y.; He, X.; Tang, F.; Dong, W. Z\*: Zero-shot Style Transfer via Attention Reweighting. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 6934–6944.
- Alami Mejjati, Y.; Richardt, C.; Tompkin, J.; Cosker, D.; Kim, K.I. Unsupervised attention-guided image-to-image translation.
   Advances in neural information processing systems 2018, 31.
- Fontanini, T.; Botti, F.; Bertozzi, M.; Prati, A. Avoiding Shortcuts in Unpaired Image-to-Image Translation. In Proceedings of the International Conference on Image Analysis and Processing. Springer, 2022, pp. 463–475.
- 25. Zhu, P.; Abdal, R.; Qin, Y.; Wonka, P. Sean: Image synthesis with semantic region-adaptive normalization. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 5104–5113.
- He, K.; Gkioxari, G.; Dollár, P.; Girshick, R. Mask r-cnn. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.
- Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You only look once: Unified, real-time object detection. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 779–788.
- Ji, X.; Henriques, J.F.; Vedaldi, A. Invariant Information Distillation for Unsupervised Image Segmentation and Clustering. CoRR 2018, abs/1807.06653, [1807.06653].

362

365

366

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

300

400

401

402

403

404

405

406

411

- Cho, J.H.; Mall, U.; Bala, K.; Hariharan, B. Picie: Unsupervised semantic segmentation using invariance and equivariance in 29. 419 clustering. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, 420 pp. 16794-16804. 421
- Fontanini, T.; Ferrari, C. Would your clothes look good on me? towards transferring clothing styles with adaptive instance 30. normalization. Sensors 2022, 22, 5002.
- Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; Joulin, A. Emerging properties in self-supervised vision 31. 424 transformers. In Proceedings of the Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 425 9650-9660. 426
- Krähenbühl, P.; Koltun, V. Efficient inference in fully connected crfs with gaussian edge potentials. Advances in neural information 32. 427 processing systems 2011, 24. 428
- 33. Liu, Z.; Luo, P.; Wang, X.; Tang, X. Deep Learning Face Attributes in the Wild. In Proceedings of the Proceedings of International 429 Conference on Computer Vision (ICCV), December 2015.
- 34. Caesar, H.; Uijlings, J.; Ferrari, V. Coco-stuff: Thing and stuff classes in context. In Proceedings of the Proceedings of the IEEE 431 conference on computer vision and pattern recognition, 2018, pp. 1209–1218. 432
- 35. Kingma, D.P.; Ba, J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 2014.
- 36. Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; Hochreiter, S. Gans trained by a two time-scale update rule converge to a 434 local nash equilibrium. Advances in neural information processing systems 2017, 30. 435
- 37. Zhang, R.; Isola, P.; Efros, A.A.; Shechtman, E.; Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. 436 In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 586–595. 437

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual 438 author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to 439 people or property resulting from any ideas, methods, instructions or products referred to in the content. 440

422

423

430