

Unsupervised Structure-Consistent Image-to-Image Translation

Shima Shahfar and Charalambos Poullis

Department of Computer Science and Software Engineering,
Concordia University,
Montreal, Quebec, Canada

Abstract. The Swapping Autoencoder achieved state-of-the-art performance in deep image manipulation and image-to-image translation. We improve this work by introducing a simple yet effective auxiliary module based on gradient reversal layers. The auxiliary module’s loss forces the generator to learn to reconstruct an image with an all-zero texture code, encouraging better disentanglement between the structure and texture information. The proposed attribute-based transfer method enables refined control in style transfer while preserving structural information *without* using a semantic mask. To manipulate an image, we encode both the geometry of the objects and the general style of the input images into two latent codes with an additional constraint that enforces structure consistency. Moreover, due to the auxiliary loss, training time is significantly reduced. The superiority of the proposed model is demonstrated in complex domains such as satellite images where state-of-the-art are known to fail. Lastly, we show that our model improves the quality metrics for a wide range of datasets while achieving comparable results with multi-modal image generation techniques.

Keywords: structure-consistent image-to-image translation · style transfer · training class imbalance

1 Introduction

Image-to-image translation and image manipulation techniques attracted much attention [10, 20, 24, 26, 28, 37, 38, 55, 59, 68] recently as they can have a significant effect on many different tasks. Of particular interest is creating realistic synthetic training datasets to improve models’ performance and generalization. One example that demonstrates the use of a synthetic dataset in the training of networks is presented in [66] where the authors introduce a semi-supervised approach to generate datasets for semantic segmentation.

There are a plethora of works [27, 28, 32] which report that for images containing single objects such as faces, or for images having the same semantic layout such as building facades, deep image manipulation techniques can produce realistic synthetic images. However, generating natural scenes or more visually complex images remains a challenge due to differences in the semantic layouts of the input images.

The challenge of deep image manipulation state-of-the-art with complex scenes is recognizing and learning essential features and characteristics from the input image. Structural information is typically shared or has common characteristics across different images in a dataset. On the other hand, the texture appears entangled with intrinsic image features. The standard approach to preserving the structural information is to condition the generation process on the input semantic mask using conditional image synthesis frameworks. However, that approach is not practical for image manipulation since the assumption of having access to semantic masks does not hold in most cases. Researchers explored different methods such as [37, 49], but in this work, we assume that image representations can be disentangled into the content/structure and texture/style.



Fig. 1: Our method learns structure-consistent image-to-image translation *without* requiring a semantic mask. We learn to disentangle structure and texture for applications such as style transfer and image editing tasks. The first(left) image shows the first input image, and the other images show the generated images in which the structure is retained from the first input image and the texture from the second, third, and fourth input images, respectively, shown in the inset images. Note that the tree’s structure is preserved, and its texture -in this case, the foliage’s colour and density- changes according to the texture of the second input image in the inset. Our model was not trained on any season transfer dataset.

To address this problem, we propose an auxiliary module that enforces the separation of structure from texture. This branch promotes the disentanglement of structure and texture by suppressing texture-related information in the structure code by applying a gradient reversal layer. Additionally, it encourages the emergence of deep features that are highly important for image editing tasks. Better structure preservation can also impact many applications ranging from creating a 3D synthetic simulation world, image editing, semantic image synthesis, and style transfer. More importantly, the proposed technique can remove biases from training datasets caused by class imbalances. Many benchmark datasets introduce bias [7, 40] that can limit the generalization capability of any network trained on them and significantly limit the impact of networks trained on these datasets in real-world scenarios.

This paper pursues three main objectives: 1) consistent and accurate structure preservation, 2) diverse, and 3) realistic image synthesis. Our goal is to learn multi-modal structure-consistent image-to-image translation in a fully unsupervised approach without requiring semantic segmentation masks. Our technical contributions can be summarized as follows:

- A new approach for a structure-consistent image-to-image translation that does not rely on prior knowledge on the scene geometry.
- An auxiliary module that enforces the disentanglement between the structure and texture information with an explicit loss term for penalizing the synthesis of realistic images when no texture information is provided.
- An extension of the Swapping Autoencoder model with our auxiliary module. We quantitatively and qualitatively demonstrate that our method generates synthetic images structurally consistent with the source input image.

We present experiments on several datasets, simple datasets with minimal variations in the semantic information of the training examples such as CelebAMaskHQ [34] Figure 2b, and complex datasets where the semantic information varies drastically such as the LSUN Church [60] Figure 1, and Cityscapes [7] Figure 4b. Our results demonstrate that the proposed method improves the performance at a fraction of the training time required by state-of-the-art.

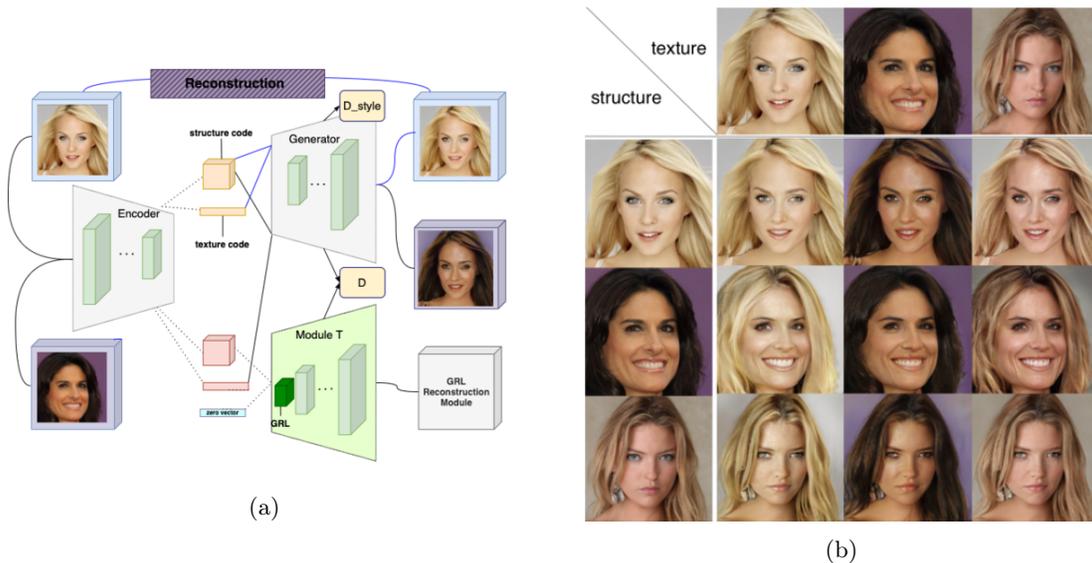


Fig. 2: (a) **Overview.** The geometry of the objects and the general style of the input images are encoded into two latent codes with an additional constraint that enforces structure consistency. We introduce a new module that encourages better disentanglement between the structure and the style, based on gradient reversal layers. This results in an attribute-based transfer that allows for a finer style transfer control while preserving structural information without requiring a semantic mask. (b) **Performance on CelebAMask-HQ:** Our model generates structure-consistent samples while transferring style from one image to another. Unlike most models that fail to preserve small structural details, our approach is able to preserve fine details such as earrings (see last row).

2 Background and Related work

This section provides an overview of the most relevant state-of-the-art, grouped according to their methodology.

Generative models. Generative Adversarial Networks (GANs) [14] introduced an adversarial process to train a generative model. The problem is formulated as a zero-sum game between a generator and discriminator where the optimal solution is to find a Nash equilibrium. Ian J. Goodfellow refers to this framework as a minimax two-player game in which generator G tries to minimize the probability of the discriminator D to recognize the fake samples, and D tries to maximize the probability of assigning the correct label. The objective function is given by,

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

GANs have proven to be very successful [4, 27, 28, 67] compared to other common approaches such as [19, 43, 47, 53, 54]. Both GANs and Variational Autoencoders (VAEs) [31] contain an encoder and a decoder; however, they differ in the sense that the GAN is a framework for estimating data distribution. On the other hand, VAEs learn the stochasticity within the data using the encoder’s latent code to match the Gaussian distribution by reparameterizing the latent distribution and maximizing the log-likelihood function. Some methods [2, 69] combine GAN and VAE or GAN and Autoencoders in their models to achieve multi-modal image generation and prevent mode collapse.

Conditional generative models such as conditional VAEs [50], conditional GANs [42], conditional autoregressive methods [15, 43], to name a few, have shown promising results [68] but we focus on conditional GANs for the rest of this section. Generative adversarial networks can be extended to conditional generative models [42] by feeding additional information c into the discriminator and generator. This c can be any information such as edge mask for semantic segmentation task or class labels for classification. By doing so, the generator can use prior noise $p_z(z)$ and additional information c to create a hidden representation and the discriminator will use the information provided as an input for a better discrimination. The quality of the results generated using conditional GANs inspired many applications employing this method, including, but not limited to, image-to-image translation [26, 38, 55, 59], image editing [5, 16], image inpainting [39, 51, 58], text-to-image [57, 63], photo colorization [36, 48, 62, 65], conditional domain adaptation [3, 5, 6, 61], super resolution [25, 33], style transfer [12, 21, 25, 27, 28, 56]. Our work extends the image-to-image translation framework with a focus on image manipulation and style transfer.

Image-to-image translation is a framework to transfer an input image into a synthesized output image while preserving some information from the input. There are many methods designed for different applications. The main difference is in the information they preserve from the input image, which depends on the application. Image-to-image translation showed promise [10, 20, 24, 68], however, as stated in [69], the quality improvement may come with the cost of losing multi-modality. Recent works show that it is possible to prevent losing multi-modality and use this method for multi-domain scenarios [22, 35, 69].

Unsupervised disentanglement aims to model the variations in data. It has been the focus of several pioneer works such as [4, 18, 49]. InfoGAN [4], for example, achieves this by maximizing the mutual information between latent variables and input data, whereas [29, 35, 45, 69] disentangle input information to structure and texture codes. Our work builds on the same principles to disentangle structure and texture in a completely unsupervised approach. However, we go one step further and aim for better disentanglement by introducing a new module to enforce better separation between the two. We show that our approach can achieve the desired disentanglement and generate realistic and diverse images while disentangling structure from style better than previous methods.

Multi-modal image synthesis overcomes the limitation of conditional GANs ignoring the latent code, also known as mode collapse. The idea behind the multi-modal image-to-image translation is to learn a conditional distribution while generating diverse images. Early works on conditional image-to-image translation were mostly focused on producing deterministic outputs [24, 38], which limits their applicability. In Section 4, we show that our method can synthesize comparable results with the current state-of-the-art [69, 70].

Style transfer also known as texture transfer, can be defined as the problem of synthesizing an image with style extracted from the source image while preserving the semantics of the content image. Recent style transfer methods [27, 28] proposed the use of conditional normalization layers such as Conditional Instance Normalization [9] and Adaptive Instance Normalization [21] as a practical approach to transfer the global style. Normalization layers used in most style transfer methods diminish semantic information. Spatially-Adaptive Normalization [44] was introduced as a way to avoid semantic-level information loss. We propose a closely related method for preserving semantic information without having access to a segmentation mask.

3 Method

Deep image manipulation requires an architecture with excellent feature extraction capabilities that allows for better disentanglement of texture from structure later on. Using an encoder, our goal is to disentangle the structure from the texture for both input images to our model. When swapping the texture or structure codes between the two randomly sampled input images

$x_1, x_2 \in \mathbb{R}^{H \times W \times 3}$, our model can synthesize an image with the same structural information as to its content reference, but having the visual appearance or texture of the style reference image. Thus, we aim to generate realistic synthesized images where the structure for the first image is preserved while transferring the style from the second image.

Our solution comprises three key modules with two discriminators namely D and D_{style} as shown in Figure 2a: an encoder E , a generator G , and a disentanglement module T which enforces better disentanglement of the structure from the style. The encoder learns how to encode visual information into two latent codes. Similar to [45], we enforce a mapping from any combination of the two latent codes to a realistic image by training an autoencoder. The generator synthesizes realistic images using the two extracted latent codes. The disentanglement module is designed to enforce the separation of the structure from the texture. We present the details of the objective function in the subsequent sections.

3.1 Encoder

The encoder E learns a mapping from the input image to two latent codes corresponding to the structure and the texture. We use a traditional autoencoder training process. We employ a reconstruction loss to measure the difference between the original image and the synthesized version with an additional non-saturating adversarial loss [14] to enforce realistic image generation, and is defined as,

$$L_{enc}(x_1, \hat{x}_1) = L_{rec}(E, G) + L_{adv}(E, G, D) = \|x_1 - G(E(x_1))\|_1 - \log(D(G(E(x_1)))) \quad (2)$$

3.2 Generator

Assuming we have already learned how to disentangle the structure from the texture, we can pass two images x_1, x_2 to the encoder and get the latent codes z_1, z_2 where $z_1 = (z_s^1, z_t^1)$ and $z_2 = (z_s^2, z_t^2)$. We assume z_s is the encoded structure and z_t is the texture of an input image and \hat{x}_1 is the reconstructed image. The generator conditioned on the latent structure code learns to map the extracted structure and texture codes to an image. The texture code will be added through weight modulation/demodulation introduced in [28]. Swapping the two texture codes before passing them to the generator is a common method to transfer style from one image to another. To ensure that the generated image is realistic, an additional non-saturating adversarial loss [14] is added, given by,

$$L_{swap}(E, G, D) = -\log(D(G(z_s^1, z_t^2))) \quad (3)$$

3.3 Structure and texture disentanglement

The latent codes must represent the structure and texture. However, this cannot be achieved in our current setting without additional constraints to encourage consistent structure and texture disentanglement. The approach used for learning consistent texture codes is to enforce all the patches sampled from the image generated in the previous step by swapping the textures to be visually similar to patches extracted from the texture reference image [45]. We achieve this using the following loss:

$$L_{style}(E, G, D_{style}) = -\log(D_{style}(C(G(z_s^1, z_t^2)), C(x^2))) \quad (4)$$

where C is a random crop of size in the range $[\frac{1}{8}, \frac{1}{4}]$. This formulation results in learning a more consistent style transfer. Experiments have shown that this term is not enough and that better disentanglement can be achieved by enforcing the structure code not to contain texture-related information. In order to enforce structure consistency, we introduce an extra module with a gradient reversal layer as its first layer followed by a generator. Gradient reversal layer act as an identity function during forward but during backward it multiplies the gradients with -1 . This

new generator has the same architecture as the original generator, but it reconstructs an image with an all-zero texture code that is theoretically impossible. Our analysis of previous works shows that structure code contains spatial information and includes style-related information. An inconsistent encoding will cause the network to generate odd samples that do not follow the algorithms and cannot be interpreted. We train this module using a reconstruction loss and a non-saturating adversarial loss [14].

$$L_{aux}(x_1, \hat{x}_1) = L_{rec}(E, T) + L_{adv}(E, T, D) = \|x_1 - T(E(x_1))\|_1 - \log(D(T(E(x_1)))) \quad (5)$$

Adding the gradient reversal layer, as shown in [11], forces the encoder to suppress any style-related information in the structure code. It also proved to be useful in cross domain disentanglement [13]. The auxiliary loss from this branch would help the encoder to disentangle structure from texture better.

3.4 Objective function

We jointly train the encoder, generators and discriminators to optimize the final objective, which is the weighted sum of previously mentioned loss functions and is given by,

$$L_{total} = \lambda_{rec}L_{enc} + \lambda_{swap}L_{swap} + \lambda_{style}L_{style} + \lambda_{aux}L_{aux} \quad (6)$$

where λ_{rec} , λ_{swap} , λ_{style} , λ_{aux} are weights that control the importance of each term. The optimal values used for each term are discussed in Section 4.

Method	LSUN Church	#iterations
StyleGAN2 [28]	57.54	48 M
Swapping [45]	52.34	14 days x 4 V100 GPUs
Ours(validation)	51.42	5 M

Table 1: Quantitative comparison of FID and training time/number of iterations on the validation set with state-of-the-art methods. Our proposed method achieves comparable performance while it converges significantly faster.

4 Experiments

Implementation details. In all reported experiments, we randomly crop and resize the input images to 256×256 resolution. We use the Adam optimizer [30] with $\beta_1 = 0.0$, $\beta_2 = 0.99$. All reported results are computed on 4 NVIDIA TESLA P100 GPUs. The discriminator D is based on StyleGAN 2 [28] and D_{style} is based on Swapping autoencoder [45]. We experimented with different hyper-parameters for λ_{rec} , λ_{swap} , λ_{style} , λ_{aux} but in this version we simply set the loss weights to be all 1.0.

Datasets. We evaluate our method on four benchmark datasets curated for scene understanding and semantic segmentation.

- CelebAMask-HQ [34] has 30,000 face images collected from the CelebA [40] dataset. CelebAMask-HQ contains annotations for 19 classes. However, we do not use masks in our training pipeline.
- LSUN church [60] is a subset of the Large-scale Scene Understanding (LSUN) dataset. The training set contains 126,227 images. It is a challenging dataset if no preprocessing is applied due to the diversity of the images.
- Cityscapes [7] is a street view dataset collected from 50 cities across Germany. The training set contains 3000 images with fine annotations, and the test set contains 500 images. It is considered a challenging dataset for image-to-image translation because each scene may contain up to 30 classes.

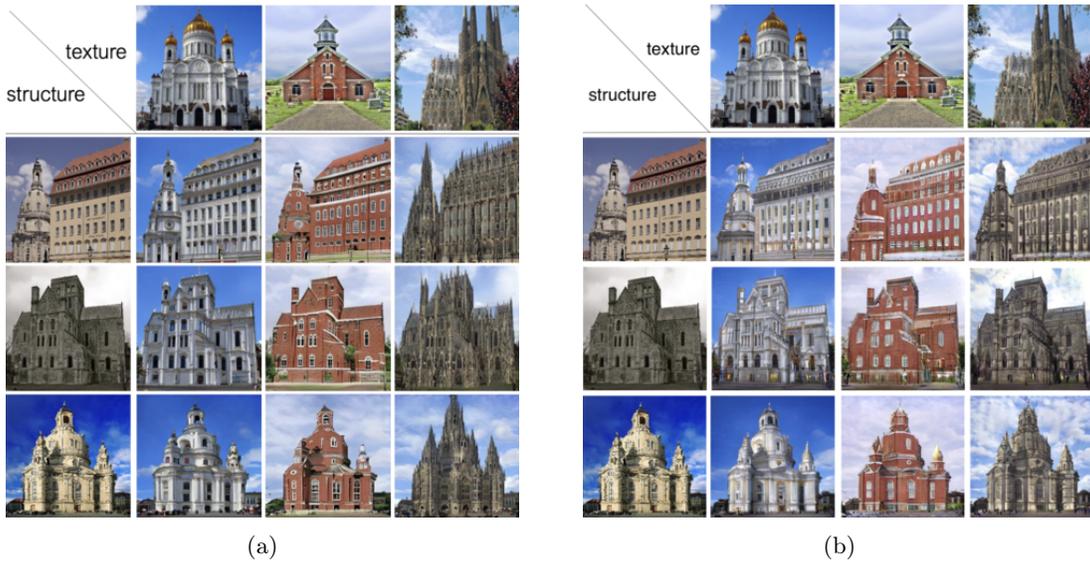


Fig. 3: Left: Results from Swapping Autoencoder [45] on LSUN Church. Right: Our results on the same images. As evident, our model achieves better feature embedding and can retain the structural information of the input image while swapping only the texture with that of a second input image. Finer-level details such as spires and buildings outline are also retained. Most notably, our model was trained for a fraction of iterations compared to [45].

- Inria [41] is an aerial imagery dataset designed for semantic segmentation of building footprints. The training set contains 180 images with 5000×5000 resolution from 5 cities. Each image covers an area of approximately $1500m \times 1500m$. The test set contains 180 images of the same size collected from 5 cities that are not part of the training set.

Baselines. We compare our approach to a number of image-to-image translation, style transfer and multi-modal image synthesis methods including Swapping Autoencoder [45], StyleGAN2 [28] and BicycleGAN [69]. We either use the results published by authors or generated using their official source code for all comparisons.

Performance metrics. We use Fréchet Inception Distance (FID) [17] to measure the quality of generated images and LPIPS [64] to compare the similarity of reconstructed images. FID calculates the difference between the real and the generated data distributions using the Inception network to extract the features while LPIPS calculates the perceptual similarity of the input with the reconstructed version. Additionally, in the Appendix(7), we report on the SIFID metric on the LSUN church dataset for the training and testing sets, and include additional comparisons and use-cases.

Structure-consistent style transfer. This section evaluates the quality of our generated images on style transfer and compares them to state-of-the-art. In Figure 3, we provide a qualitative comparison of our synthesized images with our baselines. We find that our method produces comparable results with [45] and [28] on LSUN Church dataset. A significant advantage of our approach is that it required only 5M iterations for training which demonstrates that not only is our approach significantly faster than our predecessors, but it surpasses their performance in terms of FID on the validation set, as shown in 1. Figure 3 shows that our method can generate samples with high visual quality on style transfer while preserving struc-

ture. Furthermore, structure similarity across generated samples supports the idea behind our auxiliary branch.



Fig. 4: (a) **Image translation on LSUN Church.** Each column corresponds to a particular texture extracted from the images on first row, respectively, each row contains the generated images with shared structure embedding. (b) **Image translation on Cityscapes.** The left column shows the input images from Cityscapes, the second column shows reconstruction of input images. We provide a visualization of structure latent codes in the third column after applying PCA and then resizing it to 256×256 for the purpose of visualization. The last column shows our generated images by swapping the texture between first and third row and between second and fourth row. As it can be seen, the lighting information, asphalt texture and coloring of the facades are the main information that is transferred by swapping the texture codes.

Realism of reconstruction. The diagonals of Figure 2b, 5a and 4a show the quality of our method on image reconstruction task from the learned feature embedding. Our method preserves windows, doorways, trees, spires and generally the geometry of the objects as well as finer details such as earrings and tank top strap in Figure 2b (second row). We report quantitative comparison using the LPIPS [64] to compare the similarity of reconstructed images.

Disentanglement of structure and texture. Accurately disentangling structure and texture is an important task both for style transfer and image manipulation. Given that this disentanglement is performed entirely unsupervised, we can evaluate the effectiveness of our new module by comparing the performance of our method with previous works on style transfer from existing images. Better disentanglement of structure and texture leads to a finer manipulation, resulting in significantly more realistic images. Figure 3 (left) shows the results from Swapping Autoencoder [45] on LSUN Church. Our results, shown on the right, demonstrate that our model achieves better feature embedding and generates images that retain the structural information of the input image while transferring only the texture from the second input image. Finer-level details such as spires and buildings’ outlines are also preserved.

Texture code normalization. We evaluated the effect of normalization on the texture latent code and found that applying \mathcal{L}_2 -norm results in faster convergence and more realistic synthesis. In this work we do not employ normalization in the generator, as in [23, 52], and similar to [45].

Contexts. In Figure 4a, we show examples from LSUN Church [60] that showcase the applicability of our method to other contexts. The bottom row shows a concrete example of how our technique preserves structures while transferring fine details. As it is evident, the building’s structure is preserved while the texture is replaced. Similarly, the tree’s structure is preserved, and its texture -in this case, the foliage’s colour and density- changes according to each of the source images appearing in the top row. It should be noted that the model was not trained on any season transfer dataset. Semantic image synthesis is one of the critical tasks in designing

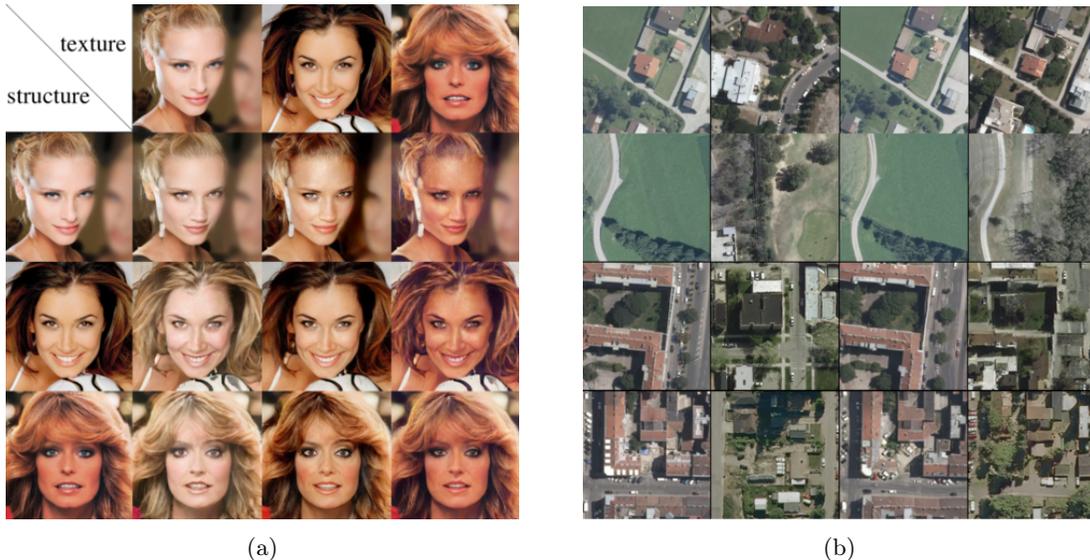


Fig. 5: (a) **Style transfer on CelebAMask-HQ.** The first row shows the texture input image. The other rows show the results using the structure image in the first column. On the second row, the specular highlight on the face is embedded as a structure and is retained. (b) **Performance on Inria dataset.** Left-to-right: first input x_1 , second input x_2 , reconstruction of x_1 , our generated sample using structure of x_1 and texture of x_2 . The semantic mask of x_1 , if available, can be transferred to the synthetic image therefore increasing the labeled images in the training set that exhibit the textural characteristics of x_2 .

3D environments, image colorization, and image editing, but it requires semantic masks and corresponding input images for training a model. This poses a limitation for many real-world applications where it is not simple to produce segmentation masks to train a conditional generative model in a supervised setting, but they need accurate semantic consistency. Our method can perfectly adopt for semantically multi-modal image synthesis in an unsupervised setting.

4.1 Comparison to state-of-the-art

Figure 9a, 9b, and 8 shows additional qualitative results on both reconstruction and style transfer tasks. The tables in Figure 6a and 6b present a quantitative comparison of our method with that of Swapping Autoencoder [45], StyleGAN2 [28], MaskGAN [34], and BicycleGAN [69].

Method	LSUN Church	CelebAMask-HQ	Cityscapes	Method	LSUN Church
Ours	51.42	29.69	162.46	StyleGAN2 [28]	0.377
Swapping [45]	52.34	32.83	182.5	Image2StyleGAN [1]	0.186
StyleGAN2 [28]	57.54	-	-	Swapping [45]	0.227
MaskGAN [34]	-	46.84	-	Ours	0.203
BicycleGAN [69]	-	-	87.74		

(a)

(b)

Fig. 6: (a) Quantitative comparison of FID on style transfer with some label-to-image translation work that are known for multimodal image synthesis and Swapping Autoencoder. In cases that we didn't have access to metric values calculated by the author, we trained their model for the same number of iterations as our network. Our method can achieve better results on CelebAMask-HQ and comparable results on LSUN Church trained for only 1.2M and 5M images. (b) Comparison of reconstructed image quality using LPIPS [64] on LSUN Church. Our method focus on preserving structural details and can produce high quality results. Given the fact that our model have only been trained on 5M images which reduce the training time by a great factor, our method can reconstruct input images better than StyleGAN2 [45].

5 Applications

As stated earlier, an important motivation of our work is to remove biases from training datasets caused by class imbalances. Benchmark datasets such as [7, 40] have inherent biases that adversely affect the network's generalization and significantly limit the effectiveness of networks used in real-world scenarios.

In this section, we present results on two unique applications employing the proposed technique:

- The first application addresses bias in training datasets and demonstrates how our method contributes to overcoming this issue.
- The second application addresses the cost-effective generation of training datasets for the task of semantic segmentation in satellite images without incurring additional labelling costs.

Furthermore, we present additional comparisons with state-of-the-art and quantitative results on the datasets LSUN Church [60], CelebAMask-HQ [34], Inria [41]. We conclude with a discussion on the limitations of our technique.

5.1 Addressing bias in training datasets

Often we talk about biases in different datasets as an issue that needs to be addressed while designing the method, and we observe some generalization issues caused mainly due to imbalances in class distributions. A different approach is to adjust or expand our existing datasets to overcome this issue. Our method can preserve fine details; for example, in face datasets, these often imbalanced features can be gender, age, skin colour, hair colour, and accessories such as earrings, eyeglasses, hats, etc. Using our method allows us to balance the dataset by generating synthetic images with under-represented features. Furthermore, in cases where labels are available for the source image, these will also be the same for the generated images since our method preserves the same structure as the source image and only changes the appearance, as shown in Figure 7.

5.2 Training datasets for Semantic Segmentation of Satellite Images

Collecting satellite imagery for semantic segmentation is known to be an expensive and challenging task. The process of capturing images is expensive, but it may also contain inaccuracies



Fig. 7: The first(left) image shows the first input image, and the second/third/fourth images show the generated image where the structure is retained from the first input image and the texture from the second/third/fourth input image, which appear in the inset images.



Fig. 8: This figure provide an example of how our method can preserve the geometry of objects and semantic details while transferring the style. This would allow us to generate multiple samples with no extra labeling cost.

due to the dynamic environment, e.g. a new building may appear that was not present at the time of acquisition of the satellite images. Another common issue is that the data collected from one city/continent cannot be easily generalized for a different city/continent. Considering all the challenges mentioned above, deploying a semantic segmentation network for aerial imagery can be challenging. Our structure-consistent network is designed to help overcome these challenges by generating realistic samples for different cities and weather conditions and generally creating datasets by style transfer. Our approach significantly reduces the time needed to process the data since we can expand any existing dataset to the desired style by only having a few images from the new city without requiring semantic labels Figure 8. Moreover, it can also be extremely useful for editing or expanding already existing datasets by changing the learned structure embedding.

6 Discussion and limitations

Our method is superior to state-of-the-art unsupervised approaches and gives comparable results to supervised techniques for image manipulation and image-to-image translation. We showed that incorporating the proposed auxiliary module as part of the training encourages better disentanglement of the structure from the texture and better feature embedding. This opens up new applications for image editing and style transfer, such as balancing existing datasets by generating images from underrepresented classes, expanding semantic segmentation datasets, creating multi-view datasets, etc. Previous works [8] explored the effect of combining multiple loss functions with different weights in a single model using [18] to achieve better optimization.

We believe the same can be applied as a future step on our pipeline for image manipulation. The importance of structure versus texture may differ from one application to another. By designing an architecture in which one can specify the percentage of structure versus texture for image generation, our method can address even broader range of challenges.

The proposed method works best when both structure and texture reference images contain the same object classes. Otherwise, the model’s behaviour is not entirely predictable. An example of this limitation is where the texture reference image does not have vegetation, but the structure reference image contains a tree. In this scenario, the network may choose to copy the original texture. Additionally, in some cases, our network will generate an image with very little change to the structure image or replace some objects due to inconsistency between represented classes in the structure and texture reference images. We have not removed such cases during training. Ignoring them can be a reasonable next step for style transfer tasks until we better understand the underlying meaning of learned texture embedding.

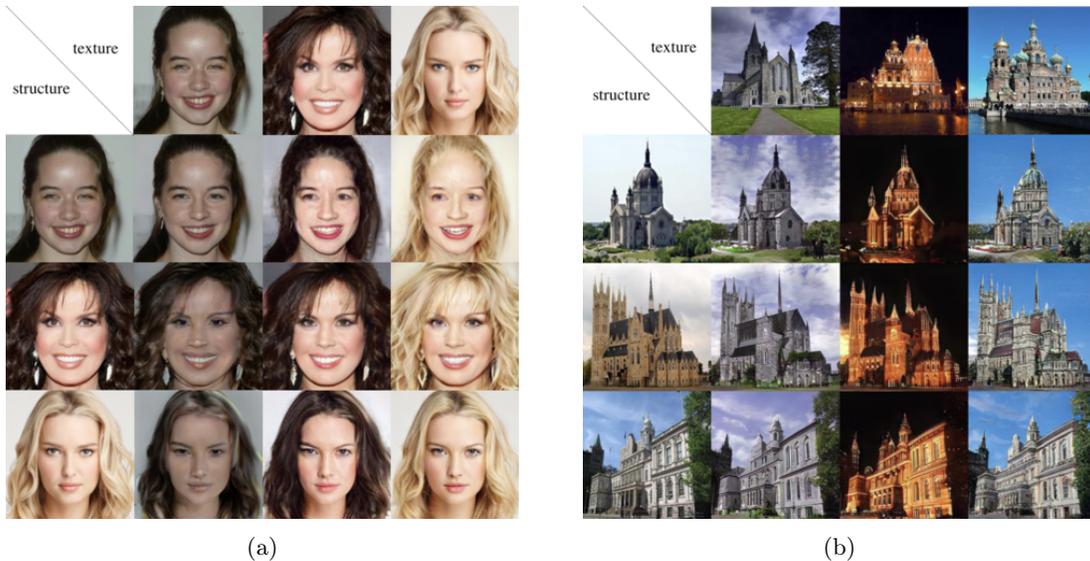


Fig. 9: (a) Examples of style transfer on CelebAMask-HQ using our learned embedding. (b) Image translation on LSUN Church showing the quality of our method in different lightning and weather.

7 Conclusions

We presented an end-to-end process for training a structure-consistent image manipulation of existing images. We showed that our approach could disentangle structure and texture with higher accuracy while preserving finer details than state-of-the-art. We have extensively tested our method and showed that it could consistently transfer texture to the correct parts and preserve structural information without requiring a semantic mask. Most notably, this is achieved while also reducing the computational time needed for training such a network to a fraction of the time needed for the current state-of-the-art. Although our method outperforms much state-of-the-art in the image-to-image translation task, defining and disentangling structure from texture in multi-object scenarios such as Cityscapes remains challenging due to the diversity of the objects and complexity of the scene. In the future, we plan to explore the knowledge embedded in latent codes for different datasets and extend this framework to other domains as discussed in Section 4.

Appendix

In this section we include additional qualitative and quantitative results, along with additional experiments, and a number of use-cases.

A. Use-case 1: Structure-guided style transfer

Image-to-image translation and image synthesis can be used for various applications requiring finer control, such as multi-view image synthesis, expanding semantic segmentation datasets to increase their variability, addressing bias in existing datasets, etc. Generating images with more control over the structure and texture can potentially address problems where collecting data is costly or where bias is present in the training set.

To verify the validity of the method and study the importance of the structure in image synthesis, we present the results in an extreme scenario where the source image used for structure does not contain texture information. To achieve this, we have used flat-shaded renders of the geometry as shown in Figure 10. The hypothesis is that given an image of the structure of the object, without any texture information, the method should preserve the fine details of the structure and transfer the related textures to the part of the image containing the structure of the object **only**. Figure 10 shows the flat-shaded render of the scene geometry without any texture as structure input and the generated images after transferring the texture information from the inset images. As shown, the proposed method successfully preserves the detailed geometry and structure while transferring the style. It should also be noted that the only area in the generated image where texture is transferred is the part of the image corresponding to the scene geometry/structure, thus experimentally proving that the structure and texture codes are disentangled. In other words, if the structure tensor was not properly separated from the texture tensor for both source images, the texture would be transferred to other areas in the image.

B. Use-case 2: Differentiable Rendering

To showcase the ability of our method ability to retain structure consistency, we perform experiments on generating multi-viewpoint imagery for differentiable rendering. Using our method, we first generate multi-view images by transferring any style to the desired structure as shown in Figures 11, 13. Next, we use a differentiable renderer to optimize the spatially varying BRDF(SVBRDF) properties from our generated multi-view images. Finally, we render using a Monte Carlo pathtracer to create renders of the geometry with the recovered SVBRDFs. Figure 12 and 14 show renders of the final results. As shown, the generated multi-view images using the proposed method preserve consistency and exhibit temporal coherence, which allows the application of differentiable rendering.



Fig. 10: **What is structure?** In this experiment we used the image on the left of size 1024x1024 from a 3D mesh with no texture to better show case what is being transferred as structure.

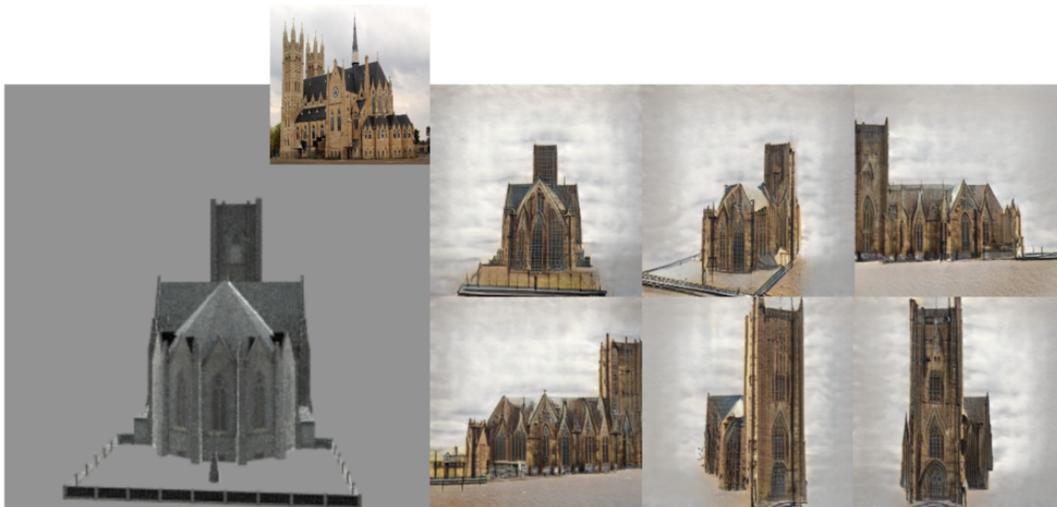


Fig. 11: Generating multi-view imagery by transferring the texture from the inset image.

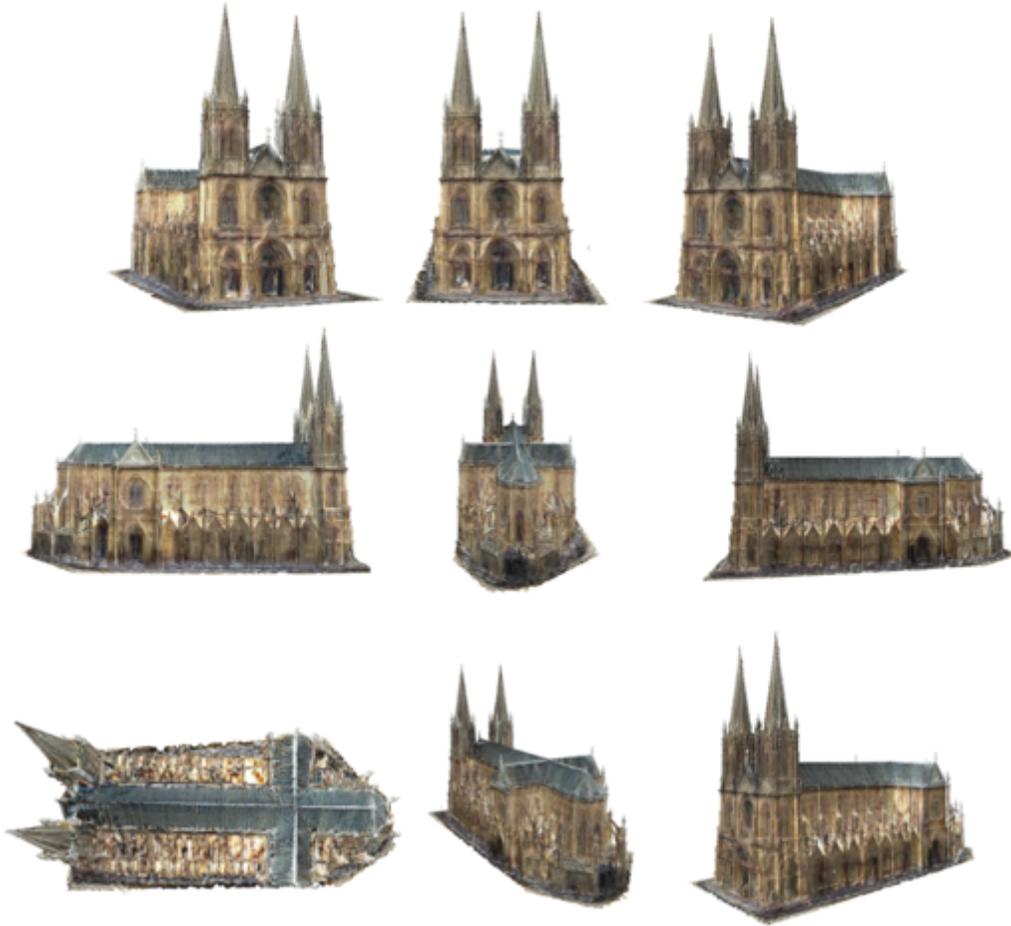


Fig. 12: Renders using the spatially varying BRDF using a differentiable renderer on our multi-view images.

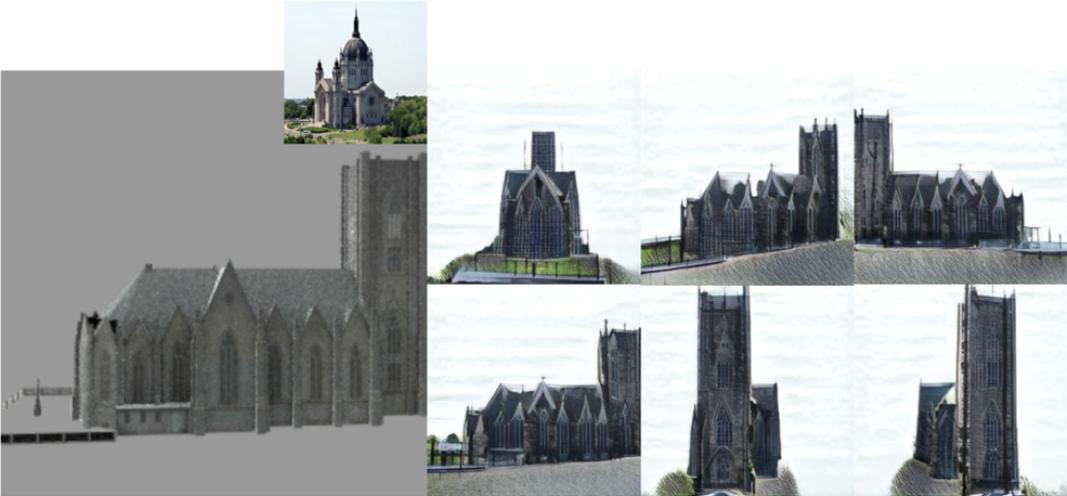


Fig. 13: Generating multi-view imagery by transferring the texture from the inset image.

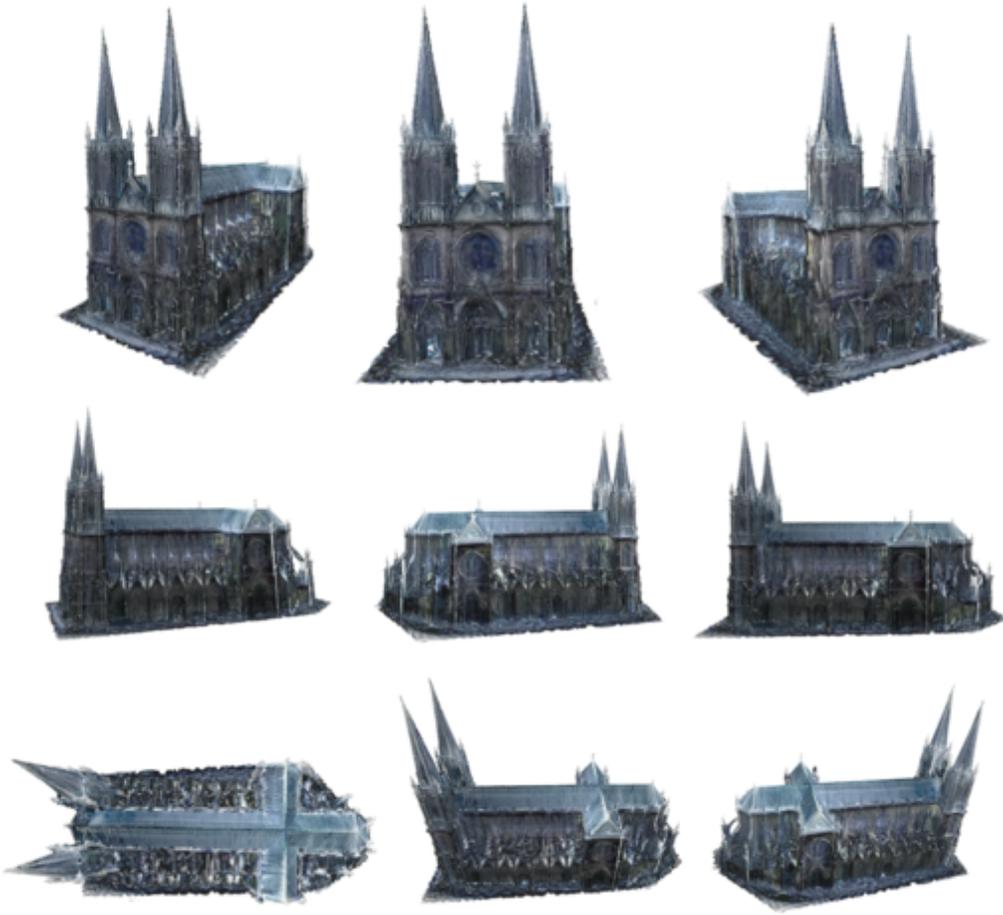


Fig. 14: Renders using the spatially varying BRDF using a differentiable renderer on our multi-view images.

C. Additional Results

Below, we present further quantitative and qualitative results.

C.1. Quantitative Results

In Figure 15 we compare the texture similarity of images generated by our proposed method with real samples based on SIFID [46].

Method	LSUN Church
StyleGAN2	52
Swapping	44
Ours	41

Fig. 15: Comparison of the performance based on **Single-Image FID** [46]

C.2. Qualitative Results

We show additional images generated by our method in Figure 16. Our method can be used to expand training datasets used for semantic segmentation by increasing the variability and reducing bias by balancing the per-class training samples, e.g. complexion, hair characteristics (color, length, style), accessories (e.g. glasses, earrings, headbands, etc.), eye color, etc. as shown in Figure 16. Similarly, Figure 17 shows the application of our method on the DeepFashion dataset. In these examples, we show some of the failures in cases where the structure is not well defined.

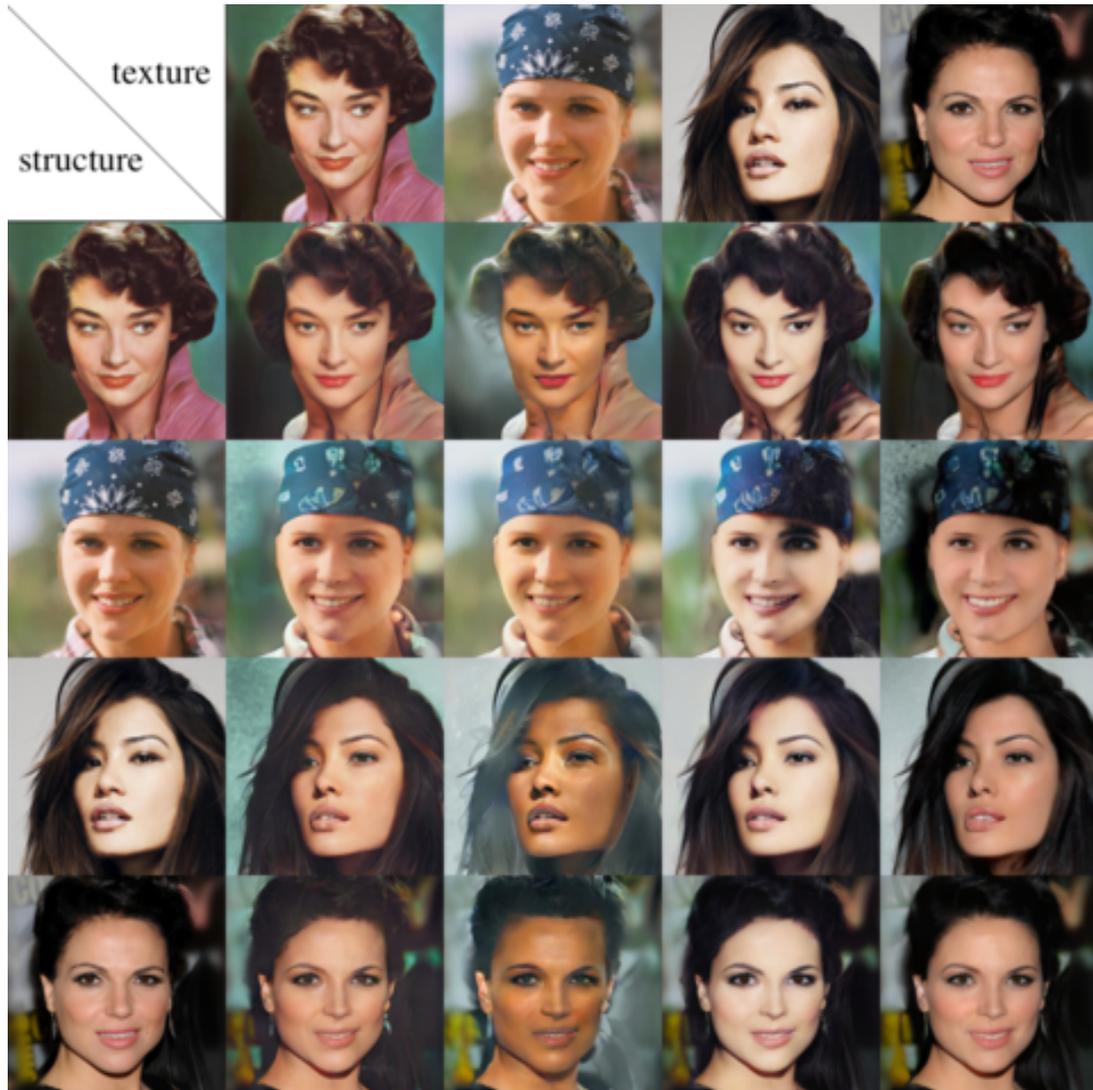


Fig. 16: Randomly selected samples from CelebAMask-HQ.



Fig. 17: Randomly selected samples from the DeepFashion dataset. Some failure cases are shown in cases where the structure is not well-defined.

References

1. Abdal, R., Qin, Y., Wonka, P.: Image2stylegan: How to embed images into the stylegan latent space? In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (October 2019) [10](#)
2. Bao, J., Chen, D., Wen, F., Li, H., Hua, G.: Cvae-gan: fine-grained image generation through asymmetric training. In: IEEE/CVF CVPR. pp. 2745–2754 (2017) [3](#)
3. Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised pixel-level domain adaptation with generative adversarial networks. In: IEEE/CVF CVPR (2017) [4](#)
4. Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., Abbeel, P.: Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In: Proceedings of the 30th International Conference on Neural Information Processing Systems. pp. 2180–2188 (2016) [3](#), [4](#)
5. 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: IEEE/CVF CVPR (June 2018) [4](#)
6. Choi, Y., Uh, Y., Yoo, J., Ha, J.W.: Stargan v2: Diverse image synthesis for multiple domains. In: IEEE/CVF CVPR (June 2020) [4](#)
7. 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: IEEE/CVF CVPR (June 2016) [2](#), [6](#), [10](#)
8. Dosovitskiy, A., Djolonga, J.: You only train once: Loss-conditional training of deep networks. In: International Conference on Learning Representations (2020) [11](#)
9. Dumoulin, V., Shlens, J., Kudlur, M.: A learned representation for artistic style. CoRR (2016) [4](#)
10. Esser, P., Haux, J., Ommer, B.: Unsupervised robust disentangling of latent characteristics for image synthesis. In: IEEE/CVF CVPR. pp. 2699–2709 (2019) [1](#), [4](#)
11. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: International conference on machine learning. pp. 1180–1189. PMLR (2015) [6](#)
12. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. 2016 IEEE CVF/CVPR pp. 2414–2423 (2016) [4](#)
13. Gonzalez-Garcia, A., Van De Weijer, J., Bengio, Y.: Image-to-image translation for cross-domain disentanglement. arXiv preprint arXiv:1805.09730 (2018) [6](#)
14. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. NeurIPS **27** (2014) [3](#), [5](#), [6](#)
15. Guadarrama, S., Dahl, R., Bieber, D., Norouzi, M., Shlens, J., Murphy, K.: Pixcolor: Pixel recursive colorization. arXiv preprint arXiv:1705.07208 (2017) [4](#)
16. He, Z., Zuo, W., Kan, M., Shan, S., Chen, X.: Attgan: Facial attribute editing by only changing what you want. IEEE TIP **28**(11), 5464–5478 (2019) [4](#)
17. 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. NeurIPS **30** (2017) [7](#)
18. Higgins, I., Matthey, L., Pal, A., Burgess, C.P., Glorot, X., Botvinick, M.M., Mohamed, S., Lerchner, A.: beta-vae: Learning basic visual concepts with a constrained variational framework. In: ICLR (2017) [4](#), [11](#)
19. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. science **313**(5786), 504–507 (2006) [3](#)
20. Hu, Q., Szabó, A., Portenier, T., Favaro, P., Zwicker, M.: Disentangling factors of variation by mixing them. In: IEEE/CVF CVPR. pp. 3399–3407 (2018) [1](#), [4](#)
21. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: ICCV (Oct 2017) [4](#)
22. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: ECCV (2018) [4](#)
23. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Bach, F., Blei, D. (eds.) Proceedings of the 32nd International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 37, pp. 448–456. PMLR, Lille, France (07–09 Jul 2015) [9](#)

24. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: *IEEE/CVF CVPR*. pp. 1125–1134 (2017) [1](#), [4](#)
25. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: *ECCV*. pp. 694–711. Springer (2016) [4](#)
26. Kaneko, T., Hiramatsu, K., Kashino, K.: Generative attribute controller with conditional filtered generative adversarial networks. *IEEE CVF/CVPR* pp. 7006–7015 (2017) [1](#), [4](#)
27. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *IEEE/CVF CVPR*. pp. 4401–4410 (2019) [1](#), [3](#), [4](#)
28. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: *Proceedings of the IEEE CVF/CVPR*. pp. 8110–8119 (2020) [1](#), [3](#), [4](#), [5](#), [6](#), [7](#), [9](#), [10](#)
29. Kazemi, H., Iranmanesh, S.M., Nasrabadi, N.: Style and content disentanglement in generative adversarial networks. In: *2019 IEEE WACV*. pp. 848–856. IEEE (2019) [4](#)
30. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. preprint arXiv:1412.6980 (2014) [6](#)
31. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013) [3](#)
32. Kotovenko, D., Sanakoyeu, A., Lang, S., Ommer, B.: Content and style disentanglement for artistic style transfer. In: *ICCV* (October 2019) [1](#)
33. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al.: Photo-realistic single image super-resolution using a generative adversarial network. In: *IEEE/CVF CVPR*. pp. 4681–4690 (2017) [4](#)
34. Lee, C.H., Liu, Z., Wu, L., Luo, P.: Maskgan: Towards diverse and interactive facial image manipulation. In: *IEEE CVF/CVPR* (2020) [2](#), [6](#), [9](#), [10](#)
35. 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) [4](#)
36. Lee, J., Kim, E., Lee, Y., Kim, D., Chang, J., Choo, J.: Reference-based sketch image colorization using augmented-self reference and dense semantic correspondence. In: *IEEE/CVF CVPR* (June 2020) [4](#)
37. Li, Y., Singh, K.K., Ojha, U., Lee, Y.J.: Mixnmatch: Multifactor disentanglement and encoding for conditional image generation. In: *IEEE/CVF CVPR* (June 2020) [1](#)
38. Liu, M.Y., Breuel, T., Kautz, J.: Unsupervised image-to-image translation networks. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *NeurIPS*. vol. 30. Curran Associates, Inc. (2017) [1](#), [4](#)
39. Liu, M.Y., Huang, X., Mallya, A., Karras, T., Aila, T., Lehtinen, J., Kautz, J.: Few-shot unsupervised image-to-image translation. In: *IEEE/CVF CVPR* (2019) [4](#)
40. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: *ICCV* (December 2015) [2](#), [6](#), [10](#)
41. Maggiori, E., Tarabalka, Y., Charpiat, G., Alliez, P.: Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark. In: *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE (2017) [7](#), [10](#)
42. Mirza, M., Osindero, S.: Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014) [4](#)
43. van den Oord, A., Kalchbrenner, N., Espeholt, L., Kavukcuoglu, K., Vinyals, O., Graves, A.: Conditional image generation with pixelcnn decoders. In: *NIPS* (2016) [3](#), [4](#)
44. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: *IEEE/CVF CVPR* (2019) [4](#)
45. Park, T., Zhu, J.Y., Wang, O., Lu, J., Shechtman, E., Efros, A.A., Zhang, R.: Swapping autoencoder for deep image manipulation. In: *NeurIPS* (2020) [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [10](#)
46. Rott Shaham, T., Dekel, T., Michaeli, T.: Singan: Learning a generative model from a single natural image. In: *Computer Vision (ICCV), IEEE International Conference on* (2019) [18](#)
47. Salimans, T., Karpathy, A., Chen, X., Kingma, D.P.: Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. arXiv preprint arXiv:1701.05517 (2017) [3](#)

48. Sangkloy, P., Lu, J., Fang, C., Yu, F., Hays, J.: Scribbler: Controlling deep image synthesis with sketch and color. In: IEEE/CVF CVPR (July 2017) 4
49. Singh, K.K., Ojha, U., Lee, Y.J.: Finegan: Unsupervised hierarchical disentanglement for fine-grained object generation and discovery. In: IEEE/CVF CVPR (2019) 1, 4
50. Sohn, K., Lee, H., Yan, X.: Learning structured output representation using deep conditional generative models. *NeurIPS* **28**, 3483–3491 (2015) 4
51. Song, Y., Yang, C., Lin, Z., Li, H., Huang, Q., Kuo, C.C.J.: Image inpainting using multi-scale feature image translation. arXiv preprint arXiv:1711.08590 (2017) 4
52. Ulyanov, D., Vedaldi, A., Lempitsky, V.S.: Instance normalization: The missing ingredient for fast stylization. *CoRR* **abs/1607.08022** (2016) 9
53. Van Oord, A., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. In: International Conference on Machine Learning. pp. 1747–1756. PMLR (2016) 3
54. Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.A.: Extracting and composing robust features with denoising autoencoders. In: Proceedings of the 25th International Conference on Machine Learning. p. 1096–1103. ICML '08, Association for Computing Machinery, New York, NY, USA (2008) 3
55. 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: Proceedings of the IEEE CVF/CVPR (2018) 1, 4
56. Xian, W., Sangkloy, P., Agrawal, V., Raj, A., Lu, J., Fang, C., Yu, F., Hays, J.: Texturegan: Controlling deep image synthesis with texture patches. In: IEEE/CVF CVPR (2018) 4
57. Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., He, X.: Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In: IEEE/CVF CVPR (June 2018) 4
58. Yeh, R., Chen, C., Lim, T.Y., Hasegawa-Johnson, M., Do, M.N.: Semantic image inpainting with perceptual and contextual losses. arXiv preprint arXiv:1607.07539 2(3) (2016) 4
59. Yi, Z., Zhang, H., Tan, P., Gong, M.: Dualgan: Unsupervised dual learning for image-to-image translation. In: ICCV (Oct 2017) 1, 4
60. Yu, F., Zhang, Y., Song, S., Seff, A., Xiao, J.: Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv preprint arXiv:1506.03365 (2015) 2, 6, 9, 10
61. Yu, X., Chen, Y., Liu, S., Li, T., Li, G.: Multi-mapping image-to-image translation via learning disentanglement. In: *NeurIPS* (2019) 4
62. Zhang, B., He, M., Liao, J., Sander, P.V., Yuan, L., Bermak, A., Chen, D.: Deep exemplar-based video colorization (June 2019) 4
63. Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., Metaxas, D.N.: Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In: ICCV (Oct 2017) 4
64. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: IEEE/CVF CVPR (June 2018) 7, 8, 10
65. Zhang, R., Zhu, J.Y., Isola, P., Geng, X., Lin, A.S., Yu, T., Efros, A.A.: Real-time user-guided image colorization with learned deep priors. *ACM TOG* **9**(4) (2017) 4
66. Zhang, Y., Ling, H., Gao, J., Yin, K., Lafleche, J.F., Barriuso, A., Torralba, A., Fidler, S.: Datasetgan: Efficient labeled data factory with minimal human effort. In: CVPR (2021) 1
67. Zhu, J.Y., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: ECCV. pp. 597–613. Springer (2016) 3
68. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017) 1, 4
69. Zhu, J.Y., Zhang, R., Pathak, D., Darrell, T., Efros, A.A., Wang, O., Shechtman, E.: Multimodal image-to-image translation by enforcing bi-cycle consistency. In: *NeurIPS*. pp. 465–476 (2017) 3, 4, 7, 9, 10
70. Zhu, Z., Xu, Z., You, A., Bai, X.: Semantically multi-modal image synthesis. In: IEEE/CVF CVPR (June 2020) 4