

# Exploiting Spatial Dimensions of Latent in GAN for Real-time Image Editing

Hyunsu Kim<sup>1,2</sup> Yunjey Choi<sup>1</sup> Junho Kim<sup>1</sup> Sungjoo Yoo<sup>2</sup> Youngjung Uh<sup>3</sup>

<sup>1</sup>NAVER AI Lab <sup>2</sup>Seoul National University <sup>3</sup>Yonsei University

## Abstract

Generative adversarial networks (GANs) synthesize realistic images from random latent vectors. Although manipulating the latent vectors controls the synthesized outputs, editing real images with GANs suffers from i) time-consuming optimization for projecting real images to the latent vectors, ii) or inaccurate embedding through an encoder. We propose StyleMapGAN: the intermediate latent space has spatial dimensions, and a spatially variant modulation replaces AdaIN. It makes the embedding through an encoder more accurate than existing optimization-based methods while maintaining the properties of GANs. Experimental results demonstrate that our method significantly outperforms state-of-the-art models in various image manipulation tasks such as local editing and image interpolation. Last but not least, conventional editing methods on GANs are still valid on our StyleMapGAN. Source code is available at <https://github.com/naver-ai/StyleMapGAN>.

## 1. Introduction

Generative adversarial networks (GANs) [16] have evolved dramatically in recent years, enabling high-fidelity image synthesis with models that are learned directly from data [6, 24, 25]. Recent studies have shown that GANs naturally learn to encode rich semantics within the latent space, thus changing the latent code leads to manipulating the corresponding attributes of the output images [22, 42, 17, 15, 43, 3, 50, 5]. However, it is still challenging to apply these manipulations to real images since the GAN lacks an inverse mapping from an image back to its corresponding latent code.

One promising approach for manipulating real images is image-to-image translation [21, 57, 9, 26, 28], where the model learns to synthesize an output image given a user's input directly. However, these methods require pre-defined tasks and heavy supervision (e.g., input-output pairs, class labels) for training and limit the user controllability at inference time. Another approach is to utilize pretrained GAN

models by directly optimizing the latent code for an individual image [1, 2, 56, 34, 36]. However, even on high-end GPUs, it requires minutes of computation for each target image, and it does not guarantee that the optimized code would be placed in the original latent space of GAN.

A more practical approach is to train an extra encoder, which learns to project an image into its corresponding latent code [31, 55, 39, 33, 40]. Although this approach enables real-time projection in a single feed-forward manner, it suffers from the low fidelity of the projected image (i.e., losing details of the target image). We attribute this limitation to the absence of spatial dimensions in the latent space. Without the spatial dimensions, an encoder compresses an image's local semantics into a vector in an entangled manner, making it difficult to reconstruct the image (e.g., vector-based or low-resolution bottleneck layer is not capable of producing high-frequency details [30, 8]).

As a solution to such problems, we propose StyleMapGAN which exploits *stylemap*, a novel representation of the latent space. Our key idea is simple. Instead of learning a vector-based latent representation, we utilize a tensor with explicit spatial dimensions. Our proposed representation benefits from its spatial dimensions, enabling GANs to easily encode the local semantics of images into the latent space. This property allows an encoder to effectively project an image into the latent space, thus providing high-fidelity and real-time projection. Our method also offers a new capability to edit specific regions of an image by manipulating the matching positions of the stylemap.

On multiple datasets, our stylemap indeed substantially enhances the projection quality compared to the traditional vector-based latent representation (§4.3). Furthermore, we show the advantage of our method over state-of-the-art methods on image projection, interpolation, and local editing (§4.4 & §4.5). Finally, we show that our method can transplant regions even when the regions are not aligned between one image and another (§4.6).

## 2. Related work

**Optimization-based editing methods** iteratively update the latent vector of pre-trained GANs to project a real image

into the latent space [56, 7, 1, 55, 20, 4]. For example, Image2StyleGAN [1] reconstructs the image by optimizing intermediate representation for each layer of StyleGAN [24]. In-DomainGAN [55] focuses not only on reconstructing the image in pixel space, but also on landing the inverted code in the semantic domain of the original latent space. Neural Collage [47] and pix2latent [20] present a hybrid optimization strategy for projecting an image into the latent space of class-conditional GANs (e.g., BigGAN [6]). On the other hand, we exploit an encoder, which makes editing two to three orders of magnitude faster than optimization methods.

**Learning-based editing methods** train an extra encoder to directly infer the latent code given a target image [31, 13, 12, 14, 40]. For example, ALI [14] and BiGAN [12] introduce a fully adversarial framework to jointly learn the generator and the inverse mapping. Several work [31, 45, 49] has been made towards combining the variational autoencoder [29] with GANs for latent projection. ALAE [40] builds an encoder to predict the intermediate latent space of StyleGAN. However, all the above methods provide limited reconstruction quality due to the lack of spatial dimensions of latent space. Swap Autoencoder [38] learns to encode an image into two components, structure code and texture code, and generate a realistic image given any swapped combination. Although it can reconstruct images fast and precisely thanks to such representation, texture code is still a vector, which makes structured texture transfer challenging. Our editing method successfully reflects not only the color and texture but also the shape of a reference image.

**Local editing methods** tackle editing specific parts [11, 3, 58, 53, 44] (e.g., nose, background) as opposed to the most GAN-based image editing methods modifying global appearance [42, 50, 38]. For example, Editing in Style [11] tries to identify each channel’s contribution of the per-layer style vectors to specific parts. Structured Noise [3] replaces the learned constant from StyleGAN with an input tensor, which is a combination of local and global codes. However, these methods [11, 3, 5] do not target real image, which performances are degraded significantly in the real image. SEAN [58] facilitates editing real images by encoding images into the per-region style codes and manipulating them, but it requires pairs of images and segmentation masks for training. Besides, the style code is still a vector, so it has the same problem as Swap Autoencoder [38].

### 3. StyleMapGAN

Our goal is to project images to a latent space accurately with an encoder in real-time and locally manipulate images on the latent space. We propose StyleMapGAN which adopts *stylemap*, an intermediate latent space with spatial dimensions, and a spatially variant modulation based on the stylemap (§3.1). Note that the *style* denotes

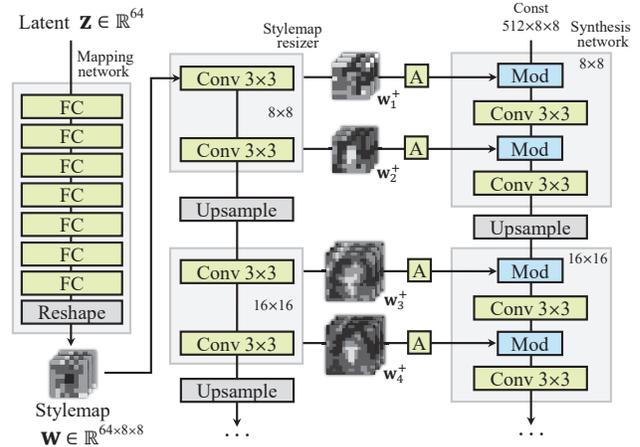


Figure 1: **StyleMapGAN Generator.** The stylemap  $w$  is resized to  $w^+$  through convolutional layers to match the spatial resolution of each feature in the synthesis network. Here “A” stands for a learned affine transform, which produces spatial modulation parameters ( $\gamma$  and  $\beta$  in Equation 1). “Mod” indicates modulation consisting of element-wise multiplication and addition. Note that the synthesis network starts from a learned constant input, and the output image’s *style* is adjusted by resized stylemaps.

not only textures (fine style) but also shapes (coarse style) following [24]. Now an encoder can embed an image to the stylemap which reconstructs the image more accurately than optimization-based methods, and partial change in the stylemap leads to local editing on the image (§3.3).

#### 3.1. Stylemap-based generator

Figure 1 describes the proposed stylemap-based generator. While a traditional mapping network produces style vectors to control feature maps, we create a stylemap with spatial dimensions, which not only makes the projection of a real image much more effective at inference but also enables local editing. The mapping network has a reshape layer at the end to produce the stylemap which forms the input to the spatially varying affine parameters. Since the feature maps in the synthesis network grow larger as getting closer to the output image, we introduce a stylemap resizer, which consists of convolutions and upsampling, to match the resolutions of stylemaps with the feature maps. The stylemap resizer resizes and transforms the stylemap with learned convolutions to convey more detailed and structured styles.

Then, the affine transform produces parameters for the modulation regarding the resized stylemaps. The modulation operation of the  $i$ -th layer in the synthesis network is as follows:

$$h_{i+1} = \left( \gamma_i \otimes \frac{h_i - \mu_i}{\sigma_i} \right) \oplus \beta_i \quad (1)$$

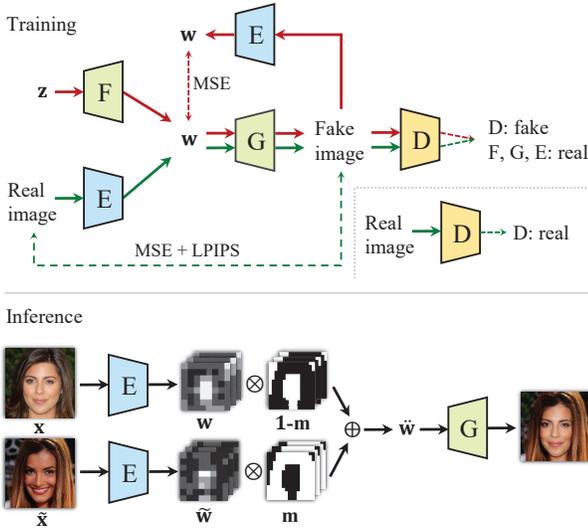


Figure 2: The upper figure contains an overall training scheme. Green and red arrows refer to flows associated with the real image and the generated image from Gaussian distribution, respectively. Dashed lines indicate loss functions. The lower figure shows our local editing method on the stylemap.

where  $\mu_i, \sigma_i \in \mathbb{R}$  are the mean and standard deviation of activations  $h_i \in \mathbb{R}^{C_i \times H_i \times W_i}$  of the  $i$ -th layer, respectively.  $\gamma_i, \beta_i \in \mathbb{R}^{C_i \times H_i \times W_i}$  are modulation parameters.  $\otimes$  and  $\oplus$  are element-wise multiplication and addition, respectively.

We remove per-pixel noise which was an extra source of spatially varying inputs in StyleGAN, because our stylemap already provides spatially varying inputs and the single input makes the projection and editing simpler. Please see the supplementary material for other details about the network and relationship with the autoencoder approach [19].

### 3.2. Training procedure and losses

In Figure 2, we use F, G, E, and D to indicate the mapping network, synthesis network with stylemap resizer, encoder, and discriminator, respectively, for brevity. D is the same as StyleGAN2, and the architecture of E is similar to D except without minibatch discrimination [41]. All networks are jointly trained using multiple losses as shown in Table 1. G and E are trained to reconstruct real images in terms of both pixel-level and perceptual-level [54]. Not only the image but E tries to reconstruct the stylemap with mean squared error (MSE) when G(F) synthesizes an image from  $z$ . D attempts to classify the real images and the fake images generated from Gaussian distribution. Lastly, we exploit domain-guided loss for the in-domain property [55]. E tries to reconstruct more realistic images by competing with D, making projected stylemap more suitable for image edit-

Loss	G	D	E
Adversarial loss [16]	✓	✓	
R <sub>1</sub> regularization [35]		✓	
Latent reconstruction			✓
Image reconstruction	✓		✓
Perceptual loss [54]	✓		✓
Domain-guided loss [55]	✓	✓	✓

Table 1: Losses for training each network. Non-saturating loss [16] is used as the adversarial loss. R1-regularization is applied every 16 steps [25] for D to stabilize training. Latent reconstruction loss is mean squared error (MSE) in the  $w$  space. Image reconstruction is MSE in image pixel-level space. We use learned perceptual image patch similarity (LPIPS) as perceptual loss for calculating perceptual differences between original and reconstructed images. Domain-guided loss is related to adversarial training that reconstructed images from the encoder tries to be classified as real by the discriminator.

ing. If we remove any of the loss functions, generation and editing performance are degraded. Refer to the supplementary material for the effect of each loss function and joint learning. Further training details are also involved.

### 3.3. Local editing

As shown at the bottom of Figure 2, the goal of local editing is to transplant some parts of a reference image to an original image with respect to a mask, which indicates the region to be modified. Note that the mask can be in any shape allowing interactive editing or label-based editing with semantic segmentation methods.

We project the original image and the reference image through the encoder to obtain stylemaps  $w$  and  $\tilde{w}$ , respectively. The edited stylemap  $\tilde{w}$  is an alpha blending of  $w$  and  $\tilde{w}$ :

$$\tilde{w} = m \otimes \tilde{w} \oplus (1 - m) \otimes w \quad (2)$$

where the mask  $m$  is shrunk by max pooling, and  $\otimes$  and  $\oplus$  are the same as Equation 1. In general, the mask is finer than  $8 \times 8$ , so we blend the stylemaps on the  $w^+$  space to achieve detailed manipulation. But for simplicity, we explain blending on the  $w$  space; the  $w^+$  space blending method is in the supplementary material. Unless otherwise stated, local editing figures are blends on the  $w^+$  space.

Contrarily to SPADE [37] or SEAN [58], even rough masks as coarse as  $8 \times 8$  produces plausible images so that the burden for user to provide detailed masks is lifted. This operation can be further revised for unidentical masks of the two images (§4.6).

## 4. Experiments

Our proposed method efficiently projects images into the style space in real-time and effectively manipulates specific regions of real images. We first describe our experimental setup (§4.1) and evaluation metrics (§4.2) and show how the proposed spatial dimensions of stylemap affect the image projection and generation quality (§4.3). We then compare our method with the state-of-the-art methods on real image projection (§4.4) and local editing (§4.5). We finally show a more flexible editing scenario and the usefulness of our proposed method (§4.6). Please see the supplementary material for high-resolution experiments and additional results such as random generation, style mixing, semantic manipulation, and failure cases.

### 4.1. Experimental setup

**Baselines.** We compare our model with recent generative models, including StyleGAN2 [25], Image2StyleGAN [1], In-DomainGAN [55], Structured Noise [3], Editing in Style [11], and SEAN [58]. We train all the baselines from scratch until they converge using the official implementations provided by the authors. For optimization-based methods [25, 1, 55, 3, 11], we use all the hyperparameters specified in their papers. We also compare our method with ALAE [40] qualitatively in the supplementary material. Note that we do not compare our method against Image2StyleGAN++ [2] and Swap Autoencoder [38], since the authors have not published their code yet.

**Datasets.** We evaluate our model on CelebA-HQ [23], AFHQ [10], and LSUN Car & Church [52]. We adopt CelebA-HQ instead of FFHQ [24], since CelebA-HQ includes segmentation masks so that we can train the SEAN baseline and exploit the masks to evaluate local editing accurately in a semantic level. The AFHQ dataset includes wider variation than the human face dataset, which is suitable for showing the generality of our model. The optimization methods take an extremely long time, we limited the test and validation set to 500 images the same as In-DomainGAN [55]. The numbers of training images for CelebA-HQ, AFHQ, and LSUN Car & Church are 29K, 15K, 5.5M, and 126K, respectively. We trained all models at  $256 \times 256$  resolution for comparison in a reasonable time, but we also provide  $1024 \times 1024$  FFHQ results in the supplementary material.

### 4.2. Evaluation metrics

**Fréchet inception distance (FID).** To evaluate the performance of image generation, we calculate FID [18] between images generated from Gaussian distribution and training set. We set the number of generated samples equal to that of training samples. We use the ImageNet-pretrained Inception-V3 [48] for feature extraction.

**FID<sub>interp</sub>.** To evaluate the global manipulation performance, we calculate FID between interpolated samples and training samples (FID<sub>interp</sub>). To generate interpolated samples, we first project 500 test images into the latent space and randomly choose pairs of latent vectors. We then generate an image using a linearly interpolated latent vector whose interpolation coefficient is randomly chosen between 0 and 1. We set the number of interpolated samples equal to that of training samples. Low FID<sub>interp</sub> indicates that the model provides high-fidelity and diverse interpolated samples.

**MSE & LPIPS.** To evaluate the projection quality, we estimate pixel-level and perceptual-level differences between target images and reconstructed images, which are mean square error (MSE) and learned perceptual image patch similarity (LPIPS) [54], respectively.

**Average precision (AP).** To evaluate the quality of locally edited images, we measure the average precision with the binary classifier trained on real and fake images [51], following the convention of the previous work [38]. We use the Blur+JPEG(0.5) model and full images for evaluation. The lower AP indicates that manipulated images are more indistinguishable from real images.

**MSE<sub>src</sub> & MSE<sub>ref</sub>.** In order to mix specific semantic, we make merged masks by combining target semantic masks of original and reference images. MSE<sub>src</sub> and MSE<sub>ref</sub> measure mean square error from the original image outside the mask and from the reference image inside the mask, respectively. To naturally combine them, images are paired by target semantic mask similarity. For local editing comparison on CelebA-HQ, 250 sets of test images are paired in each semantic (*e.g.*, background, hair) [32], which produces a total of 2500 images. For local editing on AFHQ, 250 sets of test images are paired randomly, and masks are chosen between the horizontal and vertical half-and-half mask, which produces 250 images.

### 4.3. Effects of stylemap resolution

To manipulate an image using a generative model, we first need to project the image into its latent space accurately. In Table 2, we vary the spatial resolution of stylemap and compare the performance of reconstruction and generation. For a fair comparison, we train our encoder model after training the StyleGAN2 generator. As the spatial resolution increases, the reconstruction accuracy improves significantly. It demonstrates that our stylemap with spatial dimensions is highly effective for image projection. FID varies differently across datasets, possibly due to different contextual relationships between locations for a generation. Note that our method with spatial resolution accurately preserves small details, *e.g.*, the eyes are not blurred.

Next, we evaluate the effect of the stylemap’s resolution in editing scenarios, mixing specific parts of one image and another. Figure 3 shows that the  $8 \times 8$  stylemap synthesizes



Method	Style resolution	Runtime (s)	CelebA-HQ			AFHQ		
			MSE	LPIPS	FID	MSE	LPIPS	FID
StyleGAN2	1×1	0.030	0.089	0.428	4.97	0.139	0.539	8.59
StyleMapGAN	4×4	0.085	0.062	0.351	<b>4.03</b>	0.070	0.394	14.82
StyleMapGAN	8×8	0.082	0.023	0.237	4.72	0.037	0.304	11.10
StyleMapGAN	16×16	0.078	0.010	0.146	4.71	0.016	0.183	<b>6.71</b>
StyleMapGAN	32×32	0.074	<b>0.004</b>	<b>0.076</b>	7.18	<b>0.006</b>	<b>0.090</b>	7.87

Table 2: Comparison of reconstruction and generation quality across different resolutions of the stylemap. The higher resolution helps accurate reconstruction, validating the effectiveness of stylemap. We observe that  $8 \times 8$  stylemap already provides accurate enough reconstruction and accuracy gain, and afterward, improvements get visually negligible. Although FID varies differently across datasets, possibly due to the different contextual relationships between locations for generation, the stylemap does not seriously harm the images' quality; rather, it is even better in some configurations. Using our encoder and generator, total inference time is less than 0.1s with almost perfectly reconstructed images. Although StyleGAN2 with our encoder is faster than StyleMapGAN, but it suffers from poorly reconstructed images (second column).

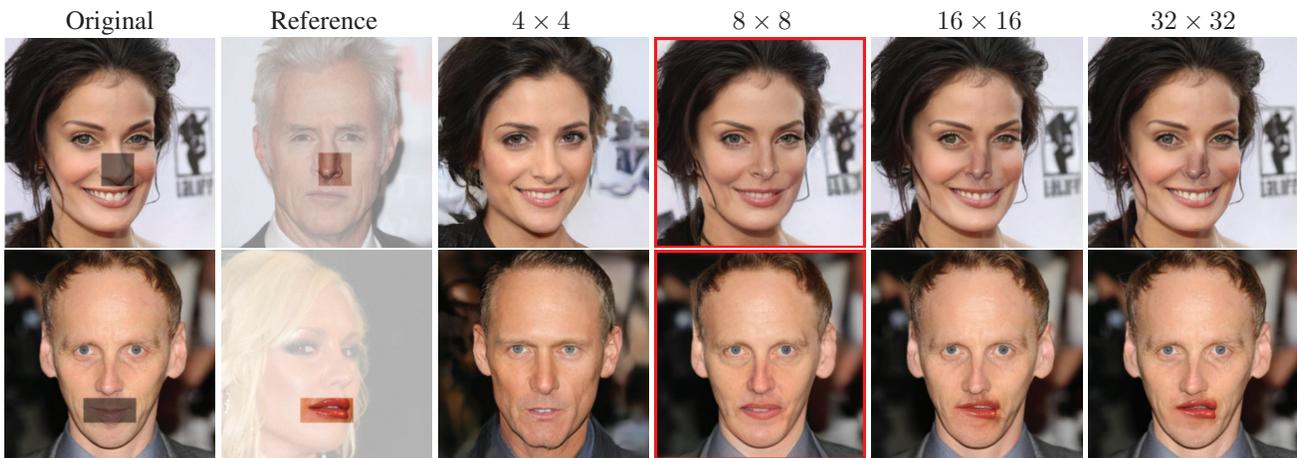
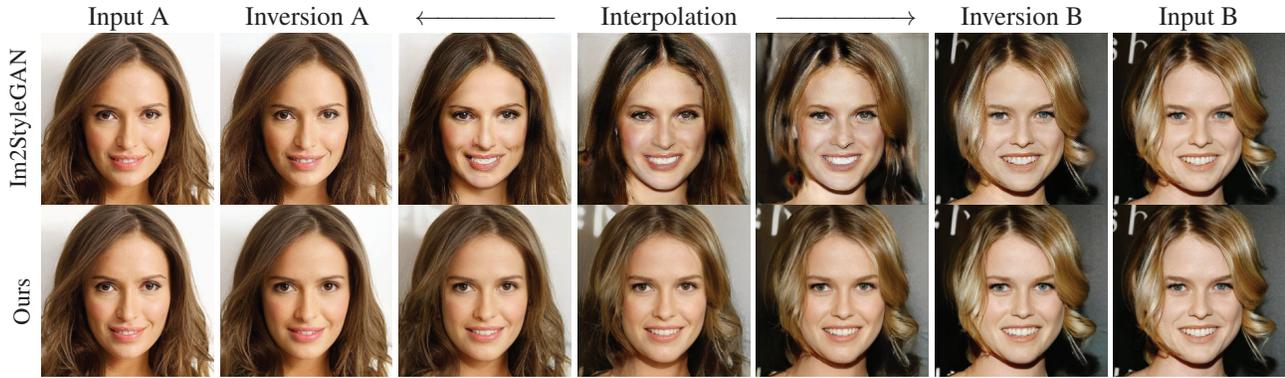


Figure 3: Local editing comparison across different resolutions of the stylemap. Regions to be discarded are faded on the original and the reference images.  $4 \times 4$  suffers from the poor reconstruction. Resolutions greater than or equal to  $16 \times 16$  result in too heterogeneous images.  $8 \times 8$  resolution shows acceptable reconstruction and natural integration. Note that our method works well even in the case that the mask locates improperly as shown in the reference image of the first row.



Method	Runtime (s)	CelebA-HQ			AFHQ		
		MSE	LPIPS	FID <sub>lerp</sub>	MSE	LPIPS	FID <sub>lerp</sub>
StyleGAN2 [25]	80.4	0.079	0.247	30.30	0.091	0.288	<b>13.87</b>
Image2StyleGAN [1]	192.5	<b>0.009</b>	<b>0.203</b>	23.68	<b>0.018</b>	<b>0.282</b>	40.80
Structured Noise [3]	64.4	0.097	0.256	27.96	0.144	0.332	34.99
In-DomainGAN [55]	6.8	0.052	0.340	<b>22.05</b>	0.077	0.414	17.54
SEAN [58]	0.146	0.064	0.334	30.29	N/A	N/A	N/A
StyleMapGAN (Ours, 8 × 8)	<b>0.082</b>	<b>0.024</b>	<b>0.242</b>	<b>9.97</b>	<b>0.037</b>	<b>0.304</b>	<b>12.42</b>

Table 3: Comparison with the baselines for real image projection. Runtime covers the end-to-end interval of projection and generation in seconds. FID<sub>lerp</sub> measures the quality of the images interpolated on the style space as a proxy for the potential quality of the manipulated images. Our method allows real-time manipulation of real images while achieving the best reconstruction accuracy and the best quality of the interpolated images. Although Image2StyleGAN produces the smallest reconstruction error, it suffers from minutes of runtime and poor interpolation quality, which are not suitable for practical editing. Its flaws can be found in the figure: rugged details in overall images, especially in teeth. SEAN is not applicable to AFHQ because it requires segmentation masks for training which are not available. The horizontal line between methods separates optimization-based methods and encoder-based methods.

the most plausible images in terms of seamlessness and preserving the identities of an original and reference image. We see that when the spatial resolution is higher than  $8 \times 8$ , the edited parts are easily detected.

Furthermore, we estimate FID<sub>lerp</sub> in different resolution models in CelebA-HQ. The  $8 \times 8$  model shows the best FID<sub>lerp</sub> value (9.97) than other resolution models; 10.72, 11.05, and 12.10 for  $4 \times 4$ ,  $16 \times 16$ , and  $32 \times 32$ , respectively. We suppose that the larger resolution of stylemap, the more likely projected latent from the encoder gets out of the latent space, which comes from a standard Gaussian distribution. Considering the editing quality and FID<sub>lerp</sub>, we choose the  $8 \times 8$  resolution as our best model and use it consistently for all subsequent experiments.

#### 4.4. Real image projection

In Table 3, we compare our approach with the state-of-the-art methods for real image projection. For both datasets, StyleMapGAN achieves better reconstruction quality (MSE & LPIPS) than all competitors except Image2StyleGAN.

However, Image2StyleGAN fails to meet requirements for editing in that it produces spurious artifacts in latent interpolation (FID<sub>lerp</sub> and figures) and suffers from minutes of runtime. Our method also achieves the best FID<sub>lerp</sub>, which implicitly shows that our manipulation on the style space leads to the most realistic images. Importantly, our method runs at least  $100\times$  faster than the optimization-based baselines since a single feed-forward pass provides accurate projection thanks to the stylemap, which is measured in a single V100 GPU. SEAN also runs with a single feed-forward pass, but it requires ground-truth segmentation masks for both training and testing, which is a severe drawback for practical uses.

#### 4.5. Local editing

We evaluate local editing performance regarding three aspects: detectability, faithfulness to the reference image in the mask, and preservation of the original image outside the mask. Figures 4 and 5 visually demonstrate that our method seamlessly composes the two images while others struggle.



Figure 4: Local editing comparison on CelebA-HQ. The first two baselines [3, 11] even fail to preserve the untouched region. In-DomainGAN loses a lot of the original image’s identity and poorly blends the two images, leaking colors to faces, hair, or background, respectively. SEAN locally transfers coarse structure and color but significantly loses details. Ours seamlessly transplants the target region from the reference to the original.

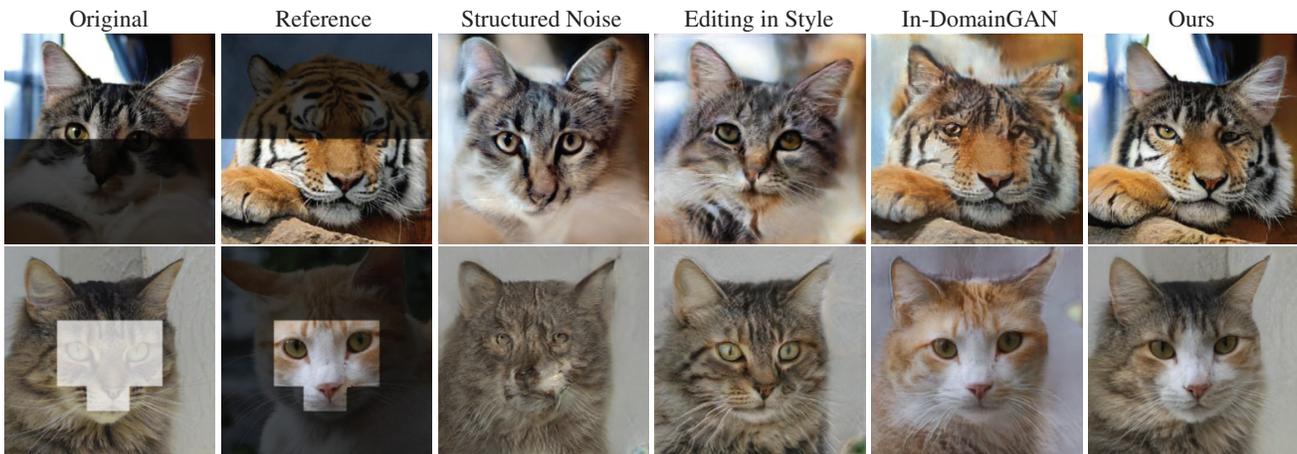


Figure 5: Local editing comparison on AFHQ. Each row blends the two images with horizontal and custom masks, respectively. Our method seamlessly composes two species with well-preserved details resulting in novel creatures, while others tend to lean towards one species.

Method	Runtime (s)	CelebA-HQ			AFHQ		
		AP	MSE <sub>src</sub>	MSE <sub>ref</sub>	AP	MSE <sub>src</sub>	MSE <sub>ref</sub>
Structured Noise [3]	64.4	99.16	0.105	0.395	99.88	0.137	0.444
Editing in Style [11]	55.6	98.34	0.094	0.321	99.52	0.130	0.417
In-DomainGAN [55]	6.8	98.72	0.164	<b>0.015</b>	99.59	0.172	<b>0.028</b>
SEAN [58]	0.155	90.41	0.067	0.141	N/A	N/A	N/A
StyleMapGAN (Ours, $8 \times 8$ )	<b>0.099</b>	<b>83.60</b>	<b>0.039</b>	<b>0.105</b>	<b>98.66</b>	<b>0.050</b>	<b>0.050</b>

Table 4: Comparison with the baselines for local image editing. Average precision (AP) is measured with the binary classifier trained on real and fake images [51]. Low AP shows our edited images are more indistinguishable from real images than other baselines. Low MSE<sub>src</sub> and MSE<sub>ref</sub> imply that our model preserves the identity of the original image and brings the characteristics of the reference image well, respectively. Our method outperforms in all metrics except MSE<sub>ref</sub> in In-DomainGAN. In-DomainGAN uses masked optimization, which only optimizes the target mask so that the identity of the original image has a great loss as shown in Figures 4 and 5.

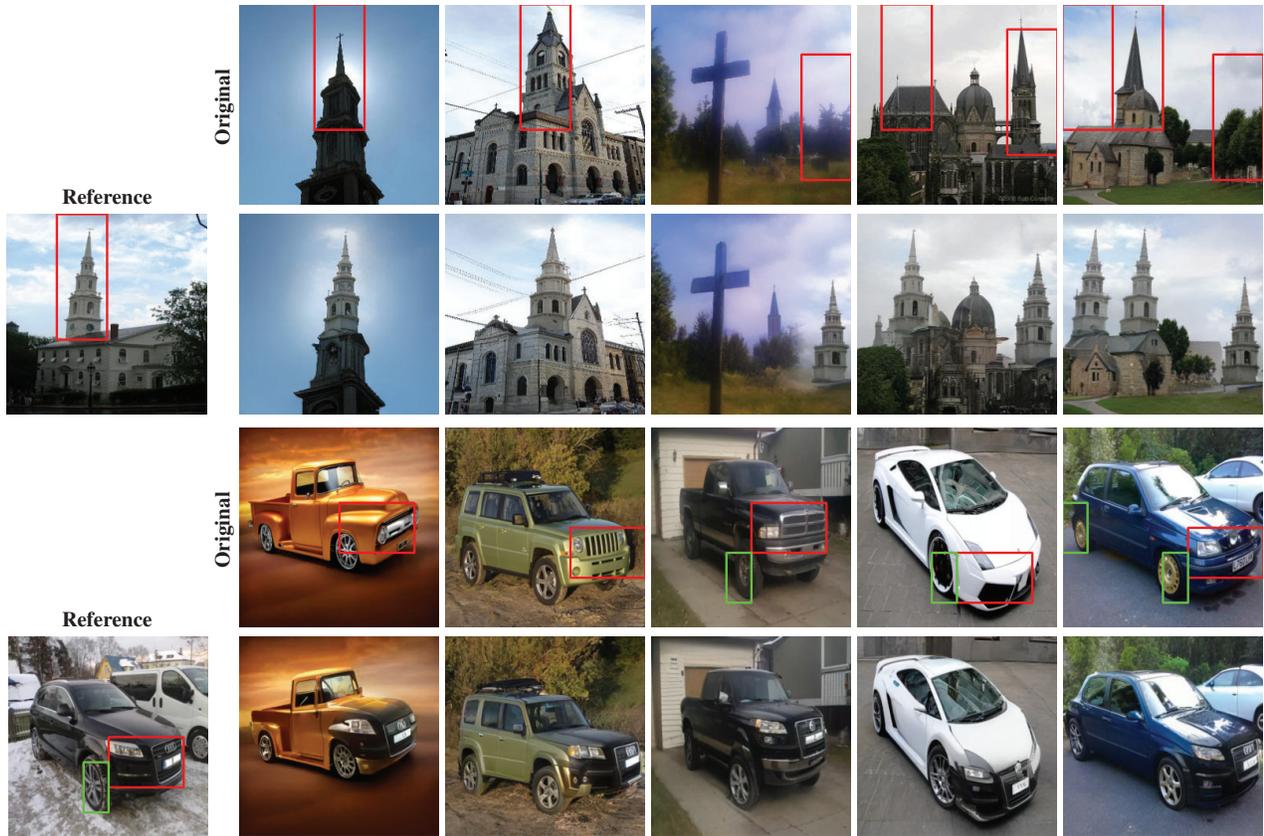


Figure 6: Examples of unaligned transplantation. StyleMapGAN allows composing arbitrary number of any regions. The size and pose of the tower, bumper and wheels are automatically adjusted regarding the surroundings. The masks are specified on  $8 \times 8$  grid and the stylemaps are blended on  $w$  space. The first row shows an example of copying one area of the reference image into multiple areas of the original images. The second row shows another example of copying two areas of the reference image. Our method can transplant the arbitrary number and size of areas of reference images.

Since there is no metric for evaluating the last two aspects, we propose two quantitative metrics:  $MSE_{src}$  and  $MSE_{ref}$ . Table 4 shows that the results from our method are the hardest for the classifier to detect fakes, and both original and reference images are best reflected. Note that MSEs are not the sole measures, but AP should be considered together for the realism of the image.

#### 4.6. Unaligned transplantation

Here, we demonstrate a more flexible use case, unaligned transplantation, showing that our local editing does not require the masks on the original and the reference images to be aligned. We project the images to the stylemaps and replace the designated region of the original stylemap with the crop of the reference stylemap even though they are in different locations. Users can specify what to replace. Figure 6 shows examples of LSUN Car & Church.

## 5. Discussion and Conclusion

Invertibility of GANs has been essential for editing real images with unconditional GAN models at a practical time, and it has not been properly answered yet. To achieve this goal, we propose StyleMapGAN, which introduces explicit spatial dimensions to the latent space, called a stylemap. We show that our method based on the stylemap has a number of advantages over prior approaches through an extensive evaluation. It can accurately project real images in real-time into the latent space and synthesize high-quality output images by both interpolation and local editing. We believe that improving fidelity by applying our latent representation to other methods such as conditional GANs (*e.g.*, BigGAN [6]) or variational autoencoders [29] would be exciting future work.

**Acknowledgements.** All experiments were conducted on NAVER Smart Machine Learning (NSML) [27, 46].

## References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *CVPR*, 2019. 1, 2, 4, 6
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *CVPR*, 2020. 1, 4
- [3] Yazeed Alharbi and Peter Wonka. Disentangled image generation through structured noise injection. In *CVPR*, 2020. 1, 2, 4, 6, 7
- [4] David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. *ACM Transactions on Graphics*, 2019. 2
- [5] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B Tenenbaum, William T Freeman, and Antonio Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. *ICLR*, 2019. 1, 2
- [6] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019. 1, 2, 8
- [7] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016. 2
- [8] Huiwen Chang, Jingwan Lu, Fisher Yu, and Adam Finkelstein. Pairedcyclegan: Asymmetric style transfer for applying and removing makeup. In *CVPR*, 2018. 1
- [9] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 1
- [10] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*, 2020. 4
- [11] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *CVPR*, 2020. 2, 4, 7
- [12] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 2
- [13] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *Advances in Neural Information Processing Systems*, pages 10542–10552, 2019. 2
- [14] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016. 2
- [15] Lore Goetschalckx, Alex Andonian, Aude Oliva, and Phillip Isola. Ganalyze: Toward visual definitions of cognitive image properties. In *ICCV*, 2019. 1
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. In *NeurIPS*, 2014. 1, 3
- [17] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. *arXiv preprint*, 2020. 1
- [18] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 4
- [19] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 2006. 3
- [20] Minyoung Huh, Richard Zhang, Jun-Yan Zhu, Sylvain Paris, and Aaron Hertzmann. Transforming and projecting images into class-conditional generative networks. *arXiv preprint arXiv:2005.01703*, 2020. 2
- [21] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial nets. In *CVPR*, 2017. 1
- [22] Ali Jahanian, Lucy Chai, and Phillip Isola. On the “steerability” of generative adversarial networks. In *ICLR*, 2020. 1
- [23] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. 2018. 4
- [24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 1, 2, 4
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 1, 3, 4, 6
- [26] Hyunsu Kim, Ho Young Jho, Eunhyeok Park, and Sungjoo Yoo. Tag2pix: Line art colorization using text tag with secat and changing loss. In *ICCV*, 2019. 1
- [27] Hanjoo Kim, Minkyu Kim, Dongjoo Seo, Jinwoong Kim, Heungseok Park, Soeun Park, Hyunwoo Jo, KyungHyun Kim, Youngil Yang, Youngkwan Kim, et al. Nsmi: Meet the mlaas platform with a real-world case study. *arXiv preprint*, 2018. 8
- [28] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwang Hee Lee. U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In *ICLR*, 2020. 1
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*, 2013. 2, 8
- [30] Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, and Marc’Aurelio Ranzato. Fader networks: Manipulating images by sliding attributes. In *NeurIPS*, 2017. 1
- [31] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, 2016. 1, 2
- [32] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 4
- [33] Junyu Luo, Yong Xu, Chenwei Tang, and Jiancheng Lv. Learning inverse mapping by autoencoder based generative adversarial nets. In *ICNIP*, 2017. 1
- [34] Fangchang Ma, Ulas Ayaz, and Sertac Karaman. Invertibility of convolutional generative networks from partial measurements. In *NeurIPS*, 2018. 1

- [35] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. Which training methods for gans do actually converge? In *ICML*, 2018. 3
- [36] Atsuhiko Noguchi and Tatsuya Harada. Image generation from small datasets via batch statistics adaptation. In *ICCV*, 2019. 1
- [37] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019. 3
- [38] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei A. Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. In *NeurIPS*, 2020. 2, 4
- [39] Guim Perarnau, Joost Van De Weijer, Bogdan Raducanu, and Jose M Álvarez. Invertible conditional gans for image editing. *arXiv preprint*, 2016. 1
- [40] Stanislav Pidhorskyi, Donald Adjeroh, and Gianfranco Doretto. Adversarial latent autoencoders. In *CVPR*, 2020. 1, 2, 4
- [41] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. 2016. 3
- [42] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *CVPR*, 2020. 1, 2
- [43] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. *arXiv preprint*, 2020. 1
- [44] Assaf Shocher, Yossi Gandelsman, Inbar Mosseri, Michal Yarom, Michal Irani, William T Freeman, and Tali Dekel. Semantic pyramid for image generation. In *CVPR*, 2020. 2
- [45] Akash Srivastava, Lazar Valkov, Chris Russell, Michael U. Gutmann, and Charles Sutton. Veegan: Reducing mode collapse in gans using implicit variational learning. In *NeurIPS*, 2017. 2
- [46] Nako Sung, Minkyu Kim, Hyunwoo Jo, Youngil Yang, Jingwoong Kim, Leonard Lausen, Youngkwan Kim, Gayoung Lee, Donghyun Kwak, Jung-Woo Ha, et al. Nsm1: A machine learning platform that enables you to focus on your models. *arXiv preprint*, 2017. 8
- [47] Ryohei Suzuki, Masanori Koyama, Takeru Miyato, Taizan Yonetsuji, and Huachun Zhu. Spatially controllable image synthesis with internal representation collaging. *arXiv preprint arXiv:1811.10153*, 2018. 2
- [48] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 4
- [49] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. It takes (only) two: Adversarial generator-encoder networks. *arXiv preprint arXiv:1704.02304*, 2017. 2
- [50] Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the gan latent space. In *ICML*, 2020. 1, 2
- [51] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *CVPR*, 2020. 4, 7
- [52] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 4
- [53] Fangneng Zhan, Hongyuan Zhu, and Shijian Lu. Spatial fusion gan for image synthesis. In *CVPR*, 2019. 2
- [54] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 3, 4
- [55] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *ECCV*, 2020. 1, 2, 3, 4, 6, 7
- [56] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *ECCV*, 2016. 1, 2
- [57] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 1
- [58] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *CVPR*, 2020. 2, 3, 4, 6, 7