

StyleSpace Analysis: Disentangled Controls for StyleGAN Image Generation

Zongze Wu

Hebrew University

zongze.wu@mail.huji.ac.il

Dani Lischinski

Hebrew University

danix@cs.huji.ac.il

Eli Shechtman

Adobe Research

elishe@adobe.com

Abstract

We explore and analyze the latent style space of StyleGAN2, a state-of-the-art architecture for image generation, using models pretrained on several different datasets. We first show that StyleSpace, the space of channel-wise style parameters, is significantly more disentangled than the other intermediate latent spaces explored by previous works. Next, we describe a method for discovering a large collection of style channels, each of which is shown to control a distinct visual attribute in a highly localized and disentangled manner. Third, we propose a simple method for identifying style channels that control a specific attribute, using a pretrained classifier or a small number of example images. Manipulation of visual attributes via these StyleSpace controls is shown to be better disentangled than via those proposed in previous works. To show this, we make use of a newly proposed Attribute Dependency metric. Finally, we demonstrate the applicability of StyleSpace controls to the manipulation of real images. Our findings pave the way to semantically meaningful and well-disentangled image manipulations via simple and intuitive interfaces.

1. Introduction

Modern Generative Adversarial Networks (GANs) are able to produce a wide variety of highly realistic synthetic images. The phenomenal success of these generative models underscores the need for a better understanding of “what makes them tick” and what kinds of control these models offer over the generated data. Of particular practical importance are controls that are interpretable and disentangled, as they suggest intuitive image manipulation interfaces.

In traditional GAN architectures, such as DCGAN [26] and Progressive GAN [16], the generator starts with a random latent vector, drawn from a simple distribution, and transforms it into a realistic image via a sequence of convolutional layers. Recently, style-based designs have become increasingly popular, where the random latent vector is first transformed into an intermediate latent code via a mapping function. This code is then used to modify the channel-

wise activation statistics at each of the generator’s convolution layers. BigGAN [6] uses class-conditional BatchNorm [14], while StyleGAN [17] uses AdaIN [13] to modulate channel-wise means and variances. StyleGAN2 [18] controls channel-wise variances by modulating the weights of the convolution kernels. It has been shown that the intermediate latent space is more disentangled than the initial one [17]. Additionally, Shen *et al.* [30] show that the latent space of StyleGAN [17, 18] is more disentangled than that of Progressive GAN [16].

Some control over the generated results may be obtained via conditioning [21], which requires training the model with annotated data. In contrast, style-based design enables discovering a variety of interpretable generator controls after training the generator. However, current methods require either a pretrained classifier [10, 30, 31, 35], a large set of paired examples [15], or manual examination of many candidate control directions [12], which limits the versatility of these approaches. Furthermore, the individual controls discovered by these methods are typically entangled, affecting multiple attributes, and are often non-local.

In this work, our goal is to understand to what degree disentanglement is inherent in style-based generator architectures. Perhaps an even more important question is to how to find these disentangled controls? In particular, can this be done in an unsupervised manner, or with only a small amount of supervision? In this paper we report several findings with respect to these questions.

Recent studies of disentangled representations [8, 28] consider a latent representation to be perfectly disentangled if each latent dimension controls a single visual attribute (*disentanglement*), and each attribute is controlled by a single dimension (*completeness*). Following this terminology, we explore the latent space of StyleGAN2 [18]. Unlike other works that analyze the (intermediate) latent space \mathcal{W} or $\mathcal{W}+$ [1], we examine *StyleSpace*, the space spanned by the channel-wise style parameters, denoted \mathcal{S} . In Section 3 we measure and compare the disentanglement and completeness of these spaces using the metrics proposed for this purpose [8]. To our knowledge we are the first to apply this quantitative framework to models trained on real

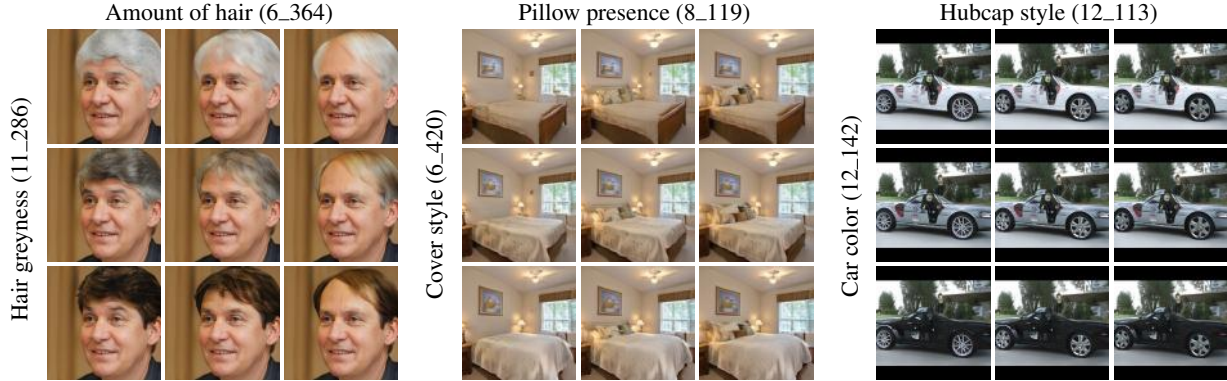


Figure 1. Disentanglement in style space, demonstrated using three different datasets. Each of the three groups above shows two manipulations that occur independently inside the same semantic region (hair, bed, and car, from left to right). The indices of the manipulated layer and channel are indicated in parentheses.

data. Our experiments reveal that \mathcal{S} is significantly better disentangled than \mathcal{W} or $\mathcal{W}+$.

In Section 4 we propose a simple method for detecting StyleSpace channels that control the appearance of local semantic regions in the image. By computing the gradient maps of generated images with respect to different style parameters, we identify those channels that are consistently active in specific semantic regions, such as hair or mouth, in the case of portraits. We demonstrate the effectiveness of this approach across three different datasets (FFHQ [17], LSUN Bedroom, and LSUN Car [37]). The StyleSpace channels that we detect are highly localized, affecting only a specific area without any visible impact of other regions. They are also surprisingly well disentangled from each other, as demonstrated in Figure 1.

Our next goal is to identify style channels that control a specific target attribute. To achieve this goal we require a set of exemplar images that exhibit the attribute of interest. The basic idea is to compare the average style vector across the exemplar set to the population average, thereby detecting dimensions that deviate the most. Our experiments indicate that such dimensions usually indeed control the target attribute, and reveal that a single attribute is typically controlled by only a few different StyleSpace channels.

To our knowledge, there is no metric to compare the disentanglement of different image manipulation controls. In Section 6 we propose Attribute Dependency (AD) as a measure for how manipulating a target attribute affects other attributes. Comparing manipulations performed in StyleSpace to those in \mathcal{W} and $\mathcal{W}+$ spaces [12, 31], shows that our controls exhibit significantly lower AD.

Finally, we share our insights about the pros and cons of two major image inversion methods, latent optimization [18, 1, 2] and encoders [39]. We show that a combination of the two may be used in order to apply our StyleSpace controls to disentangled manipulation of real images.

2. Related Work

Understanding the latent representations of pretrained generators has attracted considerable attention, since it contributes to better GAN architecture design and facilitates controllable manipulation. Bau *et al.* [5, 3] utilized semantic segmentation to analyze Progressive GAN [16] and detect causal units that control the presence of certain objects through ablation. Shen *et al.* [31] and Yang *et al.* [35] utilize classifiers to analyze StyleGAN [17] and show that a linear manipulation in \mathcal{W} space can control a specific target attribute. They further show that in $\mathcal{W}+$ space, early layers control layout, middle layers control the presence of objects, and late layers control final rendering. Collins *et al.* [7] transfer the appearance of a specific object part from a reference image to a target image, through swapping between style codes. Concurrent work by Xu *et al.* [34] shows that style space can be used for a variety of discriminative and generative tasks.

By utilizing the weights of pretrained generators, several works [1, 2, 9, 11, 23] design different latent optimization methods to do inpainting, style transfer, morphing, colorization, denoising and super resolution. Instead of latent optimization, Nitzan *et al.* [22] use the generator as a fixed decoder, and facilitate disentanglement by training an encoder for identity and another encoder for pose. Richardson *et al.* [27] do image translation by training encoders from sketches or semantic maps into StyleGAN’s \mathcal{W} space.

To facilitate attribute manipulations in an unsupervised manner, Härkönen *et al.* [12] detect interpretable controls based on PCA applied either to the latent space of StyleGAN [17] or to the feature space of BigGAN [6]. Layer-wise perturbations along the principle directions give rise to a variety of useful controls. Similarly, Shen *et al.* [32] do eigenvector decomposition in the affine transformation layer between \mathcal{W} and \mathcal{S} spaces, and use eigenvectors with

the highest eigenvalues as manipulation directions. Peebles *et al.* [24] identify interpretable controls by minimizing a Hessian loss. However, in unsupervised settings, users must examine many different manipulation directions and manually identify meaningful controls.

In contrast, we discover a large amount of localized controls using semantic maps (Section 4). The controls are ranked making it easier to detect meaningful localized manipulations in each semantic region. Furthermore, our controls are surprisingly well disentangled and fine-grained. We also detect attribute-specific controls using a small number of examples (Section 5).

3. Disentanglement of StyleGAN latent spaces

The StyleGAN/StyleGAN2 generation process involves a number of latent spaces. The first latent space, \mathcal{Z} , is typically normally distributed. Random noise vectors $z \in \mathcal{Z}$ are transformed into an *intermediate* latent space \mathcal{W} via a sequence of fully connected layers. The \mathcal{W} space is claimed to better reflect the disentangled nature of the learned distribution [17]. Each $w \in \mathcal{W}$ is further transformed to channel-wise style parameters s , using a different learned affine transformation for each layer of the generator. We refer to the space spanned by these style parameters as *StyleSpace*, or \mathcal{S} . Some works make use of another latent space, $\mathcal{W}+$, where a different intermediate latent vector w is fed to each of the generator’s layers. $\mathcal{W}+$ is mainly used for style mixing [17] and for image inversion [1, 18, 39].

In StyleGAN2 [18], there is a single style parameter per channel, which controls the feature map variances by modulating the convolution kernel weights. Additional style parameters are used by the tRGB blocks that transform feature maps to RGB images at each resolution [18]. Thus, in a 1024×1024 StyleGAN2 with 18 layers, \mathcal{W} has 512 dimensions, $\mathcal{W}+$ has 9216 dimensions, and \mathcal{S} has 9088 dimensions in total, consisting of 6048 dimensions applied to feature maps, and 3040 additional dimensions for tRGB blocks. See Appendix A for more detail. Below we refer to individual *dimensions* of \mathcal{S} as *StyleSpace channels*.

Our first goal is to determine which of these latent spaces offers the most disentangled representation. To this end, we use the recently proposed DCI (disentanglement / completeness / informativeness) metrics [8], which are suitable for comparing latent representations with different dimensions. The DCI metrics employ regressors trained using a set of latent vectors paired with corresponding attribute vectors (split into training and testing sets). *Disentanglement* measures the degree to which each latent dimension captures at most one attribute, *completeness* measures the degree to which each attribute is controlled by at most one latent dimension, while *informativeness* measures the classification accuracy of the attributes, given the latent representation.

Rather than analyzing the degree of disentanglement us-

Comparison w/ \mathcal{Z} and \mathcal{W}				Comparison with $\mathcal{W}+$			
	Disent.	Compl.	Inform.		Disent.	Compl.	Inform.
\mathcal{Z}	0.31	0.21	0.72	$\mathcal{W}+$	0.54	0.64	0.94
\mathcal{W}	0.54	0.57	0.97	\mathcal{S}	0.63	0.81	0.98
\mathcal{S}	0.75	0.87	0.99				

Table 1. Disentanglement, completeness and informativeness for different latent spaces (larger is better, maximum is 1). The two comparisons are performed using different sets of images; thus, the scores are not comparable between the two tables.

ing a synthetically generated dataset, where the factors of variations are few and known [8], we analyze StyleGAN2 trained on a real dataset, specifically FFHQ. To generate the training data for the DCI regressors, we employ 40 binary classifiers pretrained on the CelebA attributes [17]. The classifiers are trained to detect common features in portraits such as gray hair, smiling, and lipstick, and their logit outcome is converted to a binary one via a sigmoid activation.

We first randomly sample 500K latent vectors $z \in \mathcal{Z}$ and record their corresponding w and s vectors, as well as the generated images. Each image is then annotated by each of the 40 classifiers, where we record the logit, rather than just the binary outcome. Since not all attributes are well represented in the generated images (for example, there are very few portraits with a necktie), we only consider 31 attributes for which there are more than 5% positive and 5% negative outcomes. Similarly to Shen *et al.* [31], we reduce classifier uncertainty by using only the most positive 2% and most negative 2% examples, for each attribute, and split the examples equally into training and testing sets.

Finally, we compute the DCI metrics [8] to compare the latent spaces \mathcal{Z} , \mathcal{W} and \mathcal{S} . As shown in Table 1 (left), while the informativeness of both \mathcal{W} and \mathcal{S} is high and comparable, \mathcal{S} scores much higher in terms of disentanglement and completeness. This indicates that each dimension of \mathcal{S} is more likely to control a single attribute and vice versa.

Since $\mathcal{W}+$ is often used for StyleGAN inversion [1, 18], we also perform a separate experiment to compare between $\mathcal{W}+$ and \mathcal{S} . Specifically, we first randomly sample 500K intermediate latent codes $w \in \mathcal{W}$, and construct each $w+$ by concatenating n_l random w codes ($n_l = 18$ for a 1024×1024 StyleGAN2). The resulting images are somewhat less natural than those obtained in the standard manner, resulting in a smaller number of considered attributes (25 instead of 31), which we use to evaluate $\mathcal{W}+$ and \mathcal{S} as before. Table 1 (right) shows, again, that \mathcal{S} scores higher than $\mathcal{W}+$.

To our knowledge, we are the first to perform a quantitative evaluation of latent space disentanglement for a GAN model trained on real data. Since our analysis indicates that the style space \mathcal{S} is more disentangled than the other latent spaces of StyleGAN2, we proceed to further analyze \mathcal{S} in the remainder of this paper.

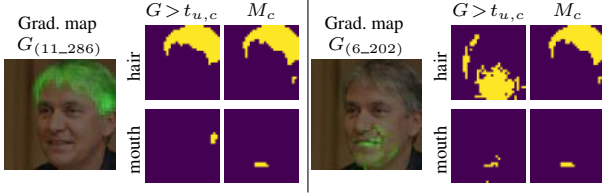


Figure 2. A gradient map with respect to each style channel u , e.g., (11_286), channel 286 of generator level 11, is thresholded against a category-specific threshold, chosen such that the resulting mask has the same size as the semantic mask M_c . The gradient mask of (11_286) has large overlap with the mask for hair, and no overlap with the mouth, while that of (6_202) has large overlap with the mask for mouth and almost none with the hair.

4. Detecting locally-active style channels

In this section we describe a simple method for detecting StyleSpace channels that control the visual appearance of local semantic regions. The intuition behind our approach is that by examining the gradient maps of generated images with respect to different channels, and measuring their overlap with specific semantic regions, we can identify those channels that are consistently active in each region. This is demonstrated in Figure 2 using two gradient maps for two different channels. If the overlap is consistent over a large number of images, these channels will be identified as locally-active for the overlapped semantic regions.

Specifically, for each image generated with style code $s \in \mathcal{S}$, we apply back-propagation to compute the gradient map of the image with respect to each channel of s . To save computation, the gradient maps are computed at a reduced spatial resolution $r \times r$ ($r = 32$ in our experiments). Next, a pretrained image segmentation network is used to obtain the semantic map M^s of the generated image. The map is resized to $r \times r$ by using the most abundant semantic category inside each bin as its semantic label. For each semantic category c and each channel u , we measure the overlap between the semantic region M_c^s and the gradient map G_u^s :

$$OC_{u,c}^s = \frac{|(G_u^s > t_{u,c}^s) \cap M_c^s|}{|M_c^s|^d}. \quad (1)$$

Here $t_{u,c}^s$ is a threshold chosen such that gradient mask $(G_u^s > t_{u,c}^s)$ has the same size as M_c^s (see Figure 2). The correction factor d gives more weight to small areas, since a large overlap between two small masks indicates precise localization. In practice, $d = 2$ gives us good balance between large and small areas.

To ensure consistency across a variety of images, we sample 1K different style codes, and compute for each code s and each channel u the semantic category with the highest overlap coefficient: $c_{s,u}^* = \arg \max_c OC_{u,c}^s$. Our goal is to detect channels for which the highest overlap category is the same for the majority of the sampled images. Furthermore,

we require that the overlap with the second most commonly affected category is twice as rare.

4.1. Experiments

We analyze StyleGAN2 [18] pretrained on FFHQ 1024x1024, LSUN Car 512x384, and LSUN Bedroom 128x128 [37]. To obtain semantic maps, we use a BiSeNet model [36] pretrained on CelebAMask-HQ [19], and a unified parsing network [33] pretrained on Broden+ [4].

As explained in Section 3 and Appendix A, 3040 channels of \mathcal{S} are used to control the tRGB blocks. None of these channels were found to have a localized effect. Rather, these channels have a global effect on the generated image, as shown in app. Figure 10.

Among the remaining 6048 channels, 1871 were found to be locally-active (in the model trained on FFHQ). Most of the detected channels control clothes (34.9%) or hair (21%). For the model trained on LSUN bedroom, we found 421 locally-active channels, most of which control the bed region (27.6%). For StyleGAN2 pretrained on LSUN car, we found 913 locally-active channels, most of which control window (33.1%) and wheel (27.3%) regions. Most of the detected channels are spread among several middle layers, with barely any channels found in early or late layers. A detailed summary of the detected locally-active channels and their breakdown by different semantic regions is included in Appendix C.

Figures 3 and 4 demonstrate some of the localized manipulations obtained by modifying the values of the channels we detected. Surprisingly, each channel appears to only control a single attribute, and even channels affecting the same local region are well disentangled, as demonstrated in Figure 1 and app. Figure 12. Unlike most controls detected by previous methods, these SpaceStyle channels provide an extremely fine-grained level of control. For example, the four channels for the ear region (last row of Figure 3), provide separate controls for the visibility of the ear, its shape, and the presence of an earring. A variety of fine-grained controls are also detected in the Car and Bedroom models (Figure 4). It should be noted that finding such interpretable disentangled local controls is very easy with our method: out of the top 10 most localized channels for each semantic area, we observe that 4–10 dimensions control (subjectively) meaningful visual attributes. A detailed breakdown by semantic category is reported in app. Table 4.

In contrast, individual channels of \mathcal{W} or $\mathcal{W}+$ space are usually entangled, with each channel affecting multiple attributes, as predicted by the DCI-based analysis from the previous section. We attribute this to the fact that each channel of $\mathcal{W}+$ affects the style parameters of an entire generation layer (via an affine transformation), rather than those of a single feature map channel.

Hair:



Mouth:



Eyes:



Eyebrows:



Ears:



Figure 3. Examples of manipulations, each controlled by a single style channel. Each pair of images shows the result of manipulation by decreasing (-) and increasing (+) the value of the style parameter (the original image is omitted). The layer index, channel index, and the direction of change is overlaid in the bottom left corner.

5. Detecting attribute-specific channels

In this section we propose a method for identifying StyleSpace channels which control a specific target attribute, specified by a set of examples. For example, given a collection of portraits of grey-haired persons, our goal is to find individual channels that control hair greyness. In contrast to InterFaceGAN [31], where around 10K positive and 10K negative examples are required, our approach typically requires only 10–30 positive exemplars. This is an important advantage, since for many attributes, negative examples can be highly varied. For example, while it is easy to find positive examples for blond hair, negative examples should ideally include all non-blond hair colors.

Our approach is based on the simple idea that the differences between the mean style vector of the positive examples (exemplar mean) and that of the entire generated distribution (population mean) reveal which channels are the most relevant for the target attribute.

Specifically, let μ^p and σ^p denote the mean and the standard deviation of the style vectors over the generated distribution. Given the style vector s^e of a specific positive example, we compute its normalized difference from the population mean: $\delta^e = \frac{s^e - \mu^p}{\sigma^p}$. Next, let μ^e and σ^e denote the mean and the standard deviation of the differences δ^e over the exemplar set. For each style channel u , the magnitude of the corresponding component μ_u^e indicates the extent to which u deviates from the population mean. Thus, we mea-

Cars:



Bedrooms:



Figure 4. Examples of manipulations, each controlled by a single style dimension. Each pair of images shows the result of manipulation by decreasing (-) and increasing (+) the value of the style parameter (the original image is omitted). The layer index, channel index, and the direction of change is overlaid in the bottom left corner.

sure the relevance of u with respect to the target attribute as the ratio $\theta_u = \frac{|\mu_u^c|}{\sigma_u^c}$. Due to the high disentanglement of \mathcal{S} (Section 3), a style channel u with a high θ_u value may be assumed to control the target attribute.

5.1. Experiments

We first use a large number (1K) of positive examples to verify that the simple method described above is indeed able to identify a set of attribute-specific control channels. Next, we demonstrate that as few as 10–30 positive examples are sufficient to detect most of these channels.

We first use the set of pretrained classifiers that were used in Section 3, to identify 1K highly positive examples for each of selected 26 attributes (see Appendix D and app. Table 5). For each attribute, we rank all the style channels (except the 3040 tRGB ones) by their relevance θ_u , and

manually examine the top 30 channels to verify that they indeed control the target attribute.

Our examination reveals that 16 out of the 26 attributes may be controlled by at least one single style channel (see supp. Table 5). The channels detected for each attribute and their ranks are reported in supp. Table 6. Interestingly, for well-defined visual attributes, such as gender, black hair, or gray hair, our method was able to find only one controlling channel. In contrast, for less specific attributes, especially those related to hair styles (bangs, receding hairline), we identified multiple controlling channels. We observe that these controls are not redundant, each controlling a unique hair style. The remaining 10 attributes are typically entangled (e.g., high cheekbones, young, or chubby), and thus no disentangled single-channel controls were detected for them. See Appendix D for further discussion.

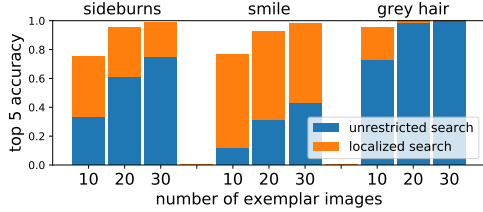


Figure 5. Top-5 detection accuracy for attribute-specific controls (for three target attributes) using 10, 20, or 30 positive examples.

Most of the detected attribute-specific control channels were highly ranked by our proposed importance score θ_u . For example, for 14 out of 16 attributes, the top-ranked channel was verified to indeed control the attribute (see app. Table 6 for the ranks of all the attribute-specific channels that we detected). This suggests that a small number of positive examples provided by a user might be sufficient for identifying such channels.

To verify the above conjecture, we randomly select sets of 10, 20, and 30 positive examples for each of three attributes (sideburns, smile, gray hair) and identify the top 5 channels for each of these small exemplar sets. If the top 5 channels include any of the verified control channels (determined using 1K images), this is considered a success. The results are reported in Figure 5.

As shown in Figure 5, increasing the number of positive examples improves the detection accuracy. The accuracy may be further improved by only considering locally-active channels (found as described in Section 4) in areas related to the target attribute. For example, if smile is the target attribute, considering only channels that are active in the mouth area, greatly improves the chances of detection. As shown by the orange bars in Figure 5, the top-5 detection rate exceeds 92% using as few as 20 examples, if the search is restricted to channels locally-active in the target area.

In summary, our approach requires only 10–30 positive examples, and detects single StyleSpace control channels. In contrast, GANSpace [12] identifies manipulation controls via a manual examination of a large number of different manipulation directions, which typically involve all of the channels of one or several layers. InterFaceGAN [31] requires more than 10K positive and 10K negative examples for each manipulation direction, which is defined in \mathcal{W} space, and thus affects all layers. Furthermore, the controls detected by these two approaches are more entangled than our control channels, as shown in the next section.

6. Disentangled attribute manipulation

In this section we compare the ability of our approach to achieve disentangled manipulation of visual attributes to that of two state-of-the-art methods, specifically GANSpace [12] and InterFaceGAN [31].

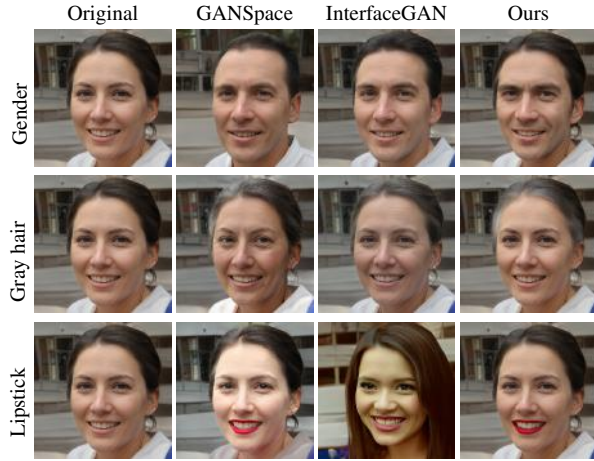


Figure 6. Comparison with state-of-the-art methods using the same amount of manipulation $\Delta l_t = 1.5\sigma(l_t)$.

Figure 6 and app. Figure 15 show a qualitative comparison between the three methods, showing the manipulation of three attributes for which the direction of manipulation is identified by all three approaches: *Gender*, *Gray hair*, and *Lipstick*. The step size along the manipulation direction is chosen such that it induces the same amount of change in the logit value l_t of the corresponding classifiers (pretrained on CelebA). Note that InterFaceGAN manipulations sometimes significantly change the identity of the person (esp. in the *Lipstick* manipulation), and some other attributes as well (added wrinkles in the *Gray hair* manipulation). GANSpace manipulations also exhibit some entanglement (*Lipstick* affects face lightness, *Gray hair* ages the rest of the face). In contrast, our approach appears to affect only the target attribute. Our *Gender* manipulation, for example, does not affect the hair style, and minimally changes the face, yet the gender unmistakably changes.

To perform a more comprehensive and quantitative comparison between the three methods, we propose a general disentanglement metric for real images, which we refer to as *Attribute Dependency* (AD). Attribute Dependency measures the degree to which manipulation along a certain direction induces changes in other attributes, as measured by classifiers for those attributes (see Appendix E for additional details). The use of classifiers here is necessary in order to cope with real images, where the exact factors of variation are not known, and we have no means to measure them. Intuitively, disentangled manipulations should induce smaller changes in other attributes.

To perform the comparison, we sample a set of images without the target attribute t (e.g., without gray hair), and manipulate them towards the target attribute, by a certain amount measured by the change in the logit outcome Δl_t of a classifier pretrained to detect attribute t . Next, we measure the change of logit Δl_i between the original images and the manipulated ones for other attributes $\forall i \in \mathcal{A} \setminus t$, where \mathcal{A} is

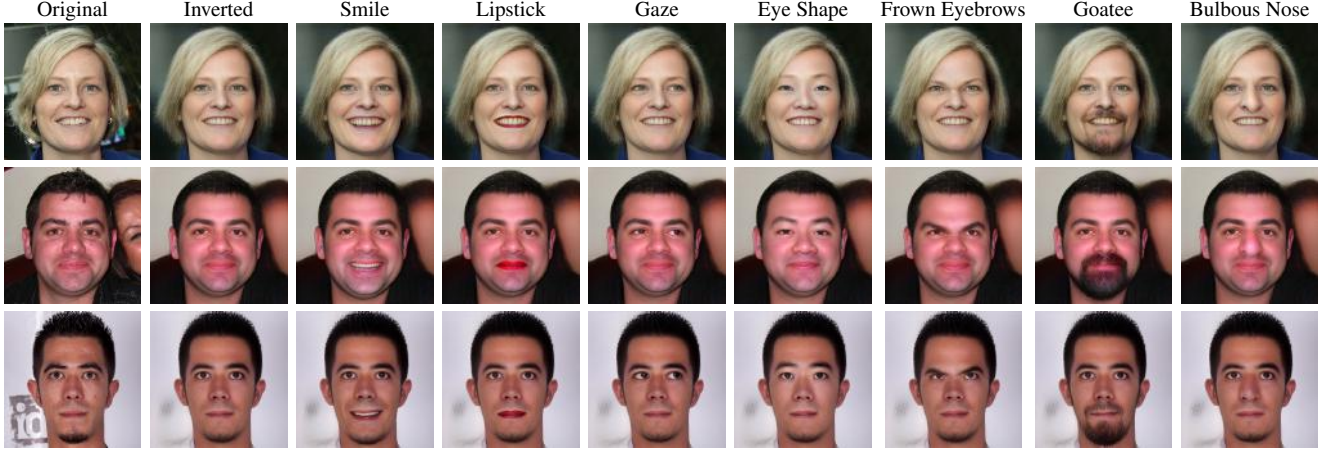


Figure 7. Manipulation of real images using encoder-based inversion. Original images are from FFHQ, and were not part of the encoder’s training set. More results can be found in app. Figure 20.

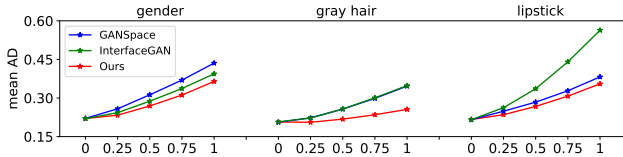


Figure 8. Mean-AD vs. the degree of target attribute manipulation ($\Delta l_t / \sigma(l_t)$). Lower mean-AD indicates better disentanglement.

the set of all attributes. Each change is normalized by $\sigma(l_i)$, the standard deviation of the logit value for attribute i over a large set of generated images. We measure mean-AD, defined as $E(\frac{1}{k} \sum_{i \in \mathcal{A} \setminus t} (\frac{\Delta l_i}{\sigma(l_i)}))$, where $k = |\mathcal{A}| - 1$. Similarly, we measure max-AD, defined as $E(\max_{i \in \mathcal{A} \setminus t} (\frac{\Delta l_i}{\sigma(l_i)}))$.

Figure 8 plots the mean-AD of the three methods (GANSpace, InterFaceGAN, and ours) for a range of manipulations of the *Gender*, *Gray hair*, and *Lipstick* attributes. It may be seen that our method (in red) exhibits a smaller mean-AD, compared to the other two methods, for each of these three attributes and across the entire manipulation range. This is consistent with our qualitative visual observations, as demonstrated in Figure 6. Our method also achieves lower max-AD scores, as reported in the supplementary material.

7. Manipulation of Real Images

To manipulate real images, it is necessary to first invert them into latent codes. This may be done via latent optimization [1, 2] or by training an encoder [40, 39] based on reconstruction loss (LPIPS [38] or L2). We adapt the latent optimization algorithm of Karras *et al.* [18] to invert real images into \mathcal{W} , $\mathcal{W}+$, and \mathcal{S} separately. Latent optimization in $\mathcal{W}+$ and \mathcal{S} spaces has more flexibility than in \mathcal{W} , enabling a closer reconstruction of the input image. Indeed, we find that the visual accuracy of the reconstruction is the highest when optimizing in \mathcal{S} , followed by $\mathcal{W}+$, and

is the lowest for \mathcal{W} (see app. Figure 18). Unfortunately, the extra flexibility may result in latent codes that do not lie on the generated image manifold, and attempting to manipulate such codes typically results in unnatural artifacts. Thus, conversely to reconstruction accuracy, we find that manipulation naturalness is best when the latent optimization is done in \mathcal{W} , followed by $\mathcal{W}+$, and the worst for \mathcal{S} (see app. Figure 19).

In order to achieve a satisfactory compromise between reconstruction accuracy and artifact-free manipulation, we train an encoder to \mathcal{S} space following the training strategy of [39] using only reconstruction loss (LPIPS). The encoder’s structure follows that of StyleALAE [25]. Due to limited computational resources, the encoder is trained on real images from FFHQ whose resolution was reduced to 128×128 . The reconstructed images bear good similarity to the input images, but exhibit some compression artifacts. The encoder’s result serves as the starting point for latent optimization [18] in \mathcal{S} space, which proceeds for a small number of iterations (50 rather than a few thousands). We find that this process can efficiently remove compression artifacts, and the resulting inversions enable artifact-free manipulation, as demonstrated in Figure 7. We believe this is because the encoder learns to embed the input real images closer to the generated image manifold, and the few optimization iterations only fine-tune the embedding.

8. Conclusion

We have shown that StyleSpace is highly disentangled, and proposed simple methods for detecting meaningful manipulation controls in this space. Future work should focus on finding meaningful control directions that involve multiple style channels. We also plan to develop inversion techniques that can deliver both high reconstruction accuracy and manipulability.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proc. ICCV*, pages 4432–4441, 2019.
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *Proc. CVPR*, pages 8296–8305, 2020.
- [3] 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 Trans. Graph.*, 38(4), 2019.
- [4] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network Dissection: quantifying interpretability of deep visual representations. In *Proc. CVPR*, pages 6541–6549, 2017.
- [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. In *Proc. ICLR*, 2019.
- [6] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [7] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of GANs. In *Proc. CVPR*, pages 5771–5780, 2020.
- [8] Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disentangled representations, 2018. In *Proc. ICLR*, volume 5, page 8, 2018.
- [9] Aviv Gabbay and Yedid Hoshen. Style generator inversion for image enhancement and animation. *arXiv preprint arXiv:1906.11880*, 2019.
- [10] Lore Goetschalckx, Alex Andonian, Aude Oliva, and Phillip Isola. GANalyze: toward visual definitions of cognitive image properties. In *Proc. ICCV*, pages 5744–5753, 2019.
- [11] Jinjin Gu, Yujun Shen, and Bolei Zhou. Image processing using multi-code GAN prior. In *Proc. CVPR*, pages 3012–3021, 2020.
- [12] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. GANSpace: discovering interpretable GAN controls. *arXiv preprint arXiv:2004.02546*, 2020.
- [13] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proc. ICCV*, pages 1501–1510, 2017.
- [14] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [15] Ali Jahanian, Lucy Chai, and Phillip Isola. On the “steerability” of generative adversarial networks. *arXiv preprint arXiv:1907.07171*, 2019.
- [16] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proc. CVPR*, pages 4401–4410, 2019.
- [18] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, pages 8110–8119, 2020.
- [19] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *Proc. CVPR*, pages 5549–5558, 2020.
- [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proc. ICCV*, pages 3730–3738, 2015.
- [21] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [22] Yotam Nitzan, Amit Bermano, Yangyan Li, and Daniel Cohen-Or. Disentangling in latent space by harnessing a pre-trained generator. *arXiv preprint arXiv:2005.07728*, 2020.
- [23] Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. *arXiv preprint arXiv:2003.13659*, 2020.
- [24] William Peebles, John Peebles, Jun-Yan Zhu, Alexei Efros, and Antonio Torralba. The hessian penalty: A weak prior for unsupervised disentanglement. *arXiv preprint arXiv:2008.10599*, 2020.
- [25] Stanislav Pidhorskyi, Donald A Adjeroh, and Gianfranco Doretto. Adversarial latent autoencoders. In *Proc. CVPR*, pages 14104–14113, 2020.
- [26] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [27] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. *arXiv preprint arXiv:2008.00951*, 2020.
- [28] Karl Ridgeway and Michael C Mozer. Learning deep disentangled embeddings with the f-statistic loss. In *Advances in Neural Information Processing Systems*, pages 185–194, 2018.
- [29] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. CVPR*, pages 815–823, 2015.
- [30] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of GANs for semantic face editing. In *Proc. CVPR*, pages 9243–9252, 2020.
- [31] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. InterFaceGAN: interpreting the disentangled face representation learned by GANs. *arXiv preprint arXiv:2005.09635*, 2020.
- [32] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in GANs. *arXiv preprint arXiv:2007.06600*, 2020.
- [33] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proc. ECCV*, pages 418–434, 2018.
- [34] Yinghao Xu, Yujun Shen, Jiapeng Zhu, Ceyuan Yang, and Bolei Zhou. Generative hierarchical features from synthesizing images. *arXiv preprint arXiv:2007.10379*, 2020.

- [35] Ceyuan Yang, Yujun Shen, and Bolei Zhou. Semantic hierarchy emerges in deep generative representations for scene synthesis. *arXiv preprint arXiv:1911.09267*, 2019.
- [36] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proc. ECCV*, pages 325–341, 2018.
- [37] 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.
- [38] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. CVPR*, pages 586–595, 2018.
- [39] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain GAN inversion for real image editing. *arXiv preprint arXiv:2004.00049*, 2020.
- [40] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *Proc. ECCV*, pages 597–613. Springer, 2016.

A. Structure of StyleGAN2 StyleSpace

To supplement the description of the different StyleGAN2 latent spaces in Section 3, here we describe the structure of the StyleSpace \mathcal{S} in more detail. Every major layer (every resolution) of the StyleGAN2 generator (synthesis network) consists of two convolution layers for feature map synthesis and a single convolution layer that converts the second feature map into an RGB image (referred to as tRGB), as shown in Figure 9. Each of these three convolution layers is modulated by a vector of style parameters. We denote the three different vectors of style parameters as s_1 , s_2 , and s_{tRGB} . These are obtained from the intermediate latent vectors $w \in \mathcal{W}$ via three affine transformations, $w_1 \rightarrow s_1$, $w_2 \rightarrow s_2$, $w_2 \rightarrow s_{tRGB}$. In \mathcal{W} space, w_1 and w_2 are the same vector, and it is the same vector for all layers. In $\mathcal{W}+$ space, w_1 and w_2 are two different vectors, and every major layer has its own pair (w_1, w_2) . The length of all the w vectors is 512. The numbers of style parameters used by the different layers are listed in Table 2. Note that in 4x4 resolution, there is only s_1 and s_{tRGB} . The length of s is 512 from the early layers until layer 14. After that layer, the length decreases from 256 to 32. In total, for a 1024x1024 generator, there are 6048 style channels that control feature maps, and 3040 additional channels that control the tRGB blocks.

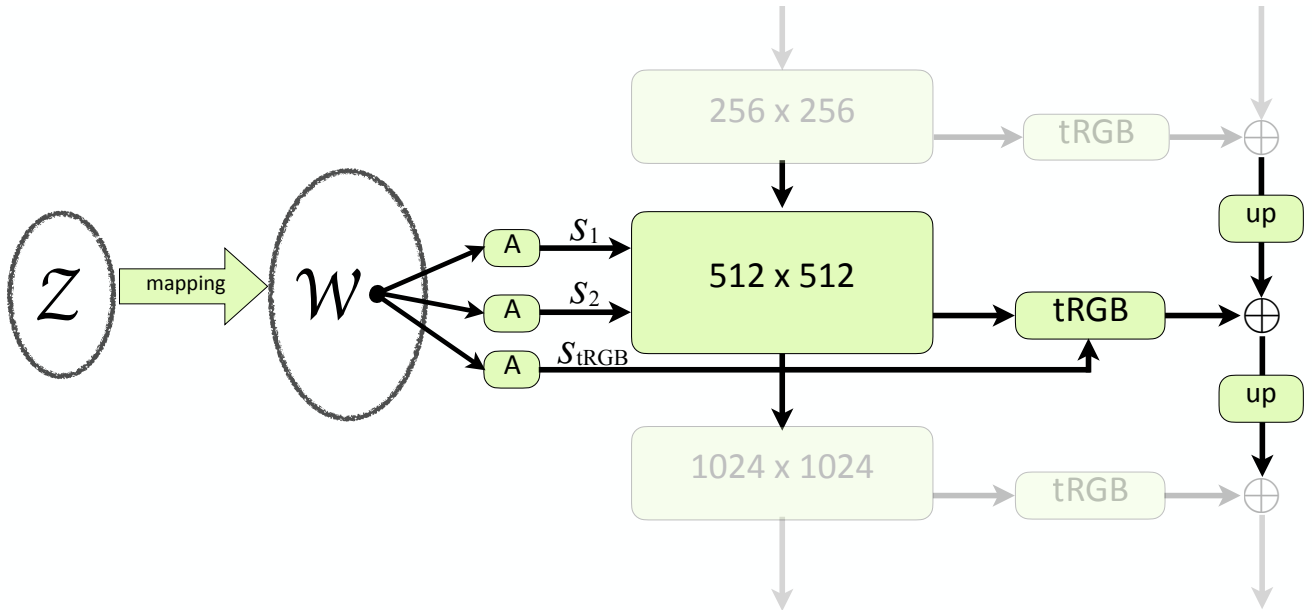


Figure 9. The internal structure of StyleSpace \mathcal{S} , shown for the 512 x 512 generator resolution.

$\mathcal{W}+$ layer index	\mathcal{S} layer index	resolution	layer name	type	# channels
0	0	4×4	Conv	s_1	512
1	1	4×4	ToRGB	s_{tRGB}	512
2	2	8×8	Conv0_up	s_1	512
3	3	8×8	Conv1	s_2	512
3	4	8×8	ToRGB	s_{tRGB}	512
4	5	16×16	Conv0_up	s_1	512
5	6	16×16	Conv1	s_2	512
5	7	16×16	ToRGB	s_{tRGB}	512
6	8	32×32	Conv0_up	s_1	512
7	9	32×32	Conv1	s_2	512
7	10	32×32	ToRGB	s_{tRGB}	512
8	11	64×64	Conv0_up	s_1	512
9	12	64×64	Conv1	s_2	512
9	13	64×64	ToRGB	s_{tRGB}	512
10	14	128×128	Conv0_up	s_1	512
11	15	128×128	Conv1	s_2	256
11	16	128×128	ToRGB	s_{tRGB}	256
12	17	256×256	Conv0_up	s_1	256
13	18	256×256	Conv1	s_2	128
13	19	256×256	ToRGB	s_{tRGB}	128
14	20	512×512	Conv0_up	s_1	128
15	21	512×512	Conv1	s_2	64
15	22	512×512	ToRGB	s_{tRGB}	64
16	23	1024×1024	Conv0_up	s_1	64
17	24	1024×1024	Conv1	s_2	32
17	25	1024×1024	ToRGB	s_{tRGB}	32

Table 2. Breakdown of StyleSpace channels by generator layers.

B. Effect of style parameters in tRGB layers

To examine the function of style parameters that control the tRGB layers, we randomly generate a set of 500K style vectors $s \in \mathcal{S}$, and perturb their s_{tRGB} channels to manipulate the tRGB layers, $s_{new} = s_{original} + n\sigma(s)$. $\sigma(s)$ is the standard deviation of each channel of s over the generated set, used to normalize the amount of perturbation across different channels [12]. n is a vector of Gaussian noise, with mean 0 and standard deviation $\sigma(n)$, which indicates the manipulation strength. Below, we use $\sigma(n) = 15$.

As shown in Figure 10, manipulating the early (coarse) resolutions (0,1,2) mainly affects the center of the target object (better visible in faces than in cars), manipulating the middle resolutions (3,4,5) typically affects the entire target object, and manipulating the late (fine) resolution layers (6,7,8) affects the entire image. (The LSUN Car model reaches only up to 512×512 , thus the fine resolution layers are (6,7)). The effect of the late (fine) resolution layers on the image is significantly stronger than that of the early and middle layers. The manipulations only affect color, without modifying shape or specific face or car related attributes.

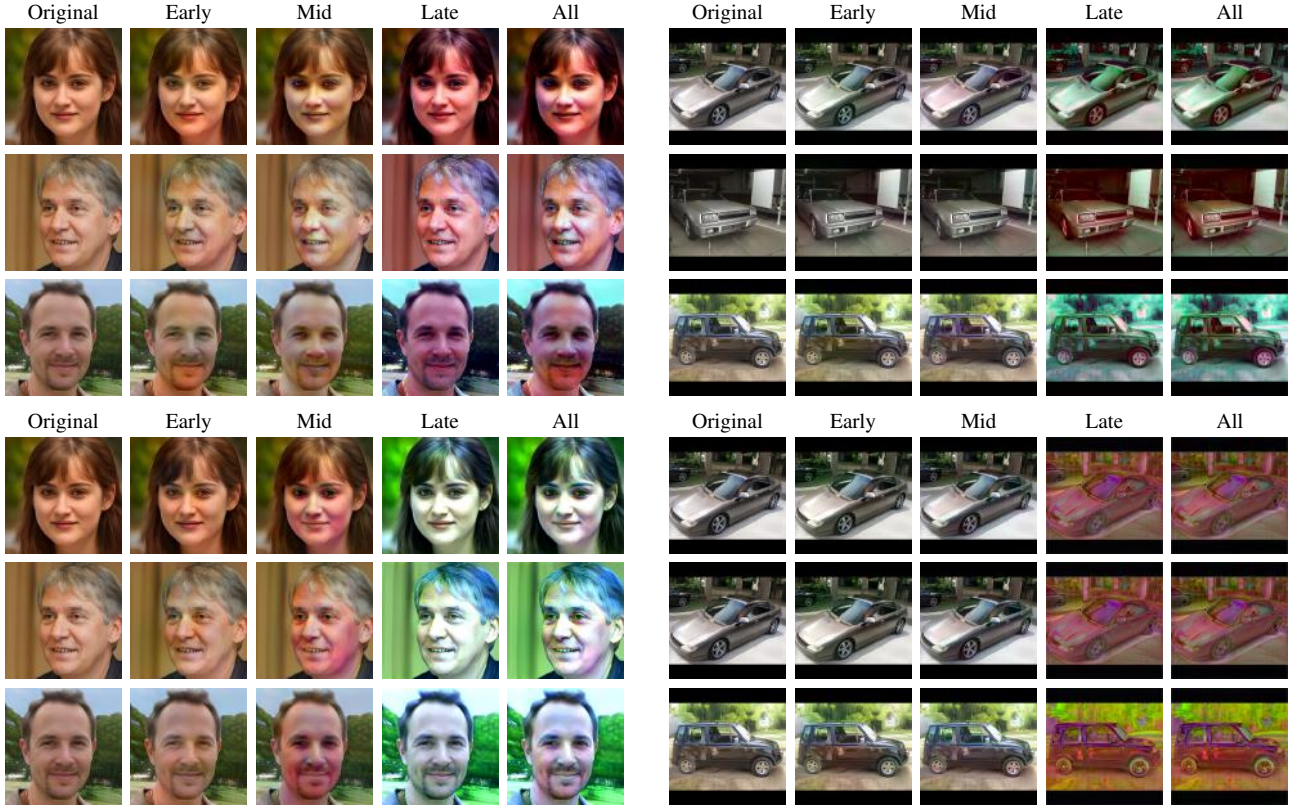


Figure 10. Manipulation by perturbing the s_{tRGB} channels in early resolution layers (0,1,2), middle layers (3,4,5) and late layers (6,7,8). Each of the four blocks above uses the same noise vector n .

C. Locally-active style channels

The number of locally-active channels that we found using the method in Section 4 for each of the three models we experimented with is summarized in Table 3. The breakdown of these localized controls across different semantic regions is plotted in Figure 11. Not all of the detected controls correspond to semantically meaningful manipulations. While there is no way to objectively determine which manipulations are meaningful, in Table 4 we report the number of manipulations that were (subjectively) determined as meaningful by the authors, among the most highly localized controls.

	FFHQ	LSUN Bedroom	LSUN Car
Num. locally-active channels	1871	421	913
Total num. of feature map style channels	6048	5376	5952
Percent of locally-active channels	30.9%	7.8%	15.3%

Table 3. Number of locally-active channels detected in different StyleGAN2 models.

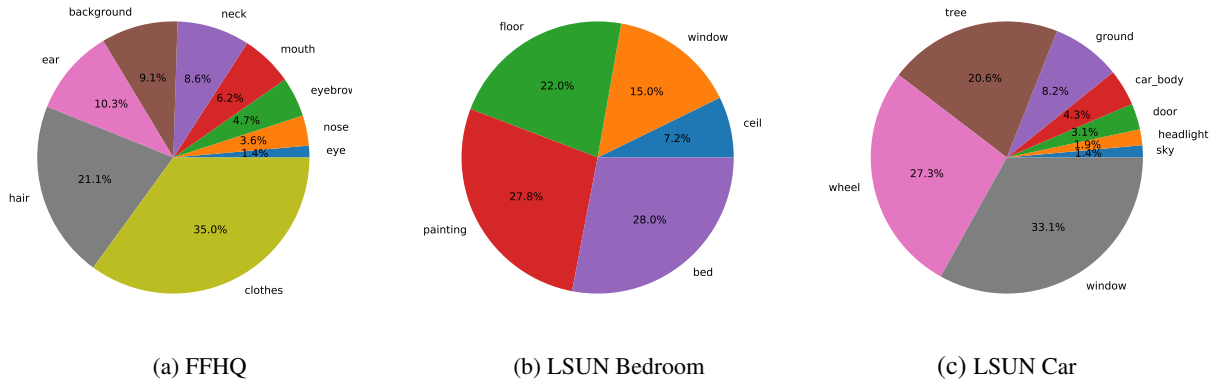


Figure 11. Breakdown of detected localized controls across different semantic regions.

	Top 5	Top 10	Top 20
Eyebrows	5	10	19
Hair	5	9	17
Nose	4	7	13
Mouth	4	7	11
Clothes	5	6	9
Neck	2	4	7
Eye	4	5	6
Ear	3	4	6
Background	5	10	15

Table 4. Number of meaningful controls among the top $k = 5, 10, 20$ most locally-active channels (those with the highest overlap coefficient, as defined by equation (1) in the main paper) in each semantic area (for the FFHQ model). Note that this count is subjective and may contain channels that control similar things (for example, size of lips).

Finally, Figure 12 demonstrates (in addition to Figure 1) the high degree of disentanglement of the localized controls that our method detects. Even pairs of controls that affect the same semantic region, typically do so in an independent manner.

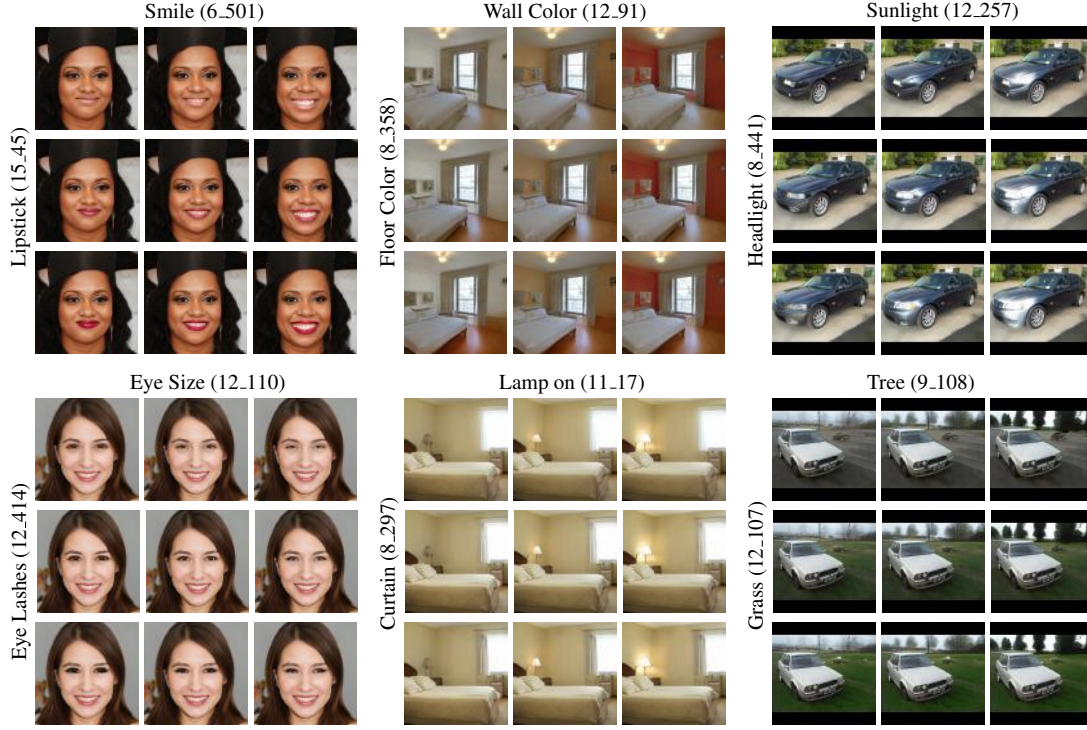


Figure 12. Disentanglement in style space, demonstrated using three different datasets (FFHQ, LSUN Bedroom, LSUN Car). Each of the six groups above shows two manipulations that occur independently inside the same image. The indices of the manipulated layer and channel are indicated in parentheses.

D. Attribute-specific channels

Starting from the 40 attributes from CelebA [17], we first remove inactivated, ambiguous and neutral attributes. Inactivated attributes (defined in Section 3) are those that are not well represented in the generated image distribution. Ambiguous attributes are highly subjective. Neutral attributes are those in between more extreme states, or attributes that are highly common across the dataset. For example, “mouth slightly open” is between an open mouth and a closed one. Another example of a neutral attribute is “no beard”, since most of the faces in FFHQ don’t have a beard. The attributes that were found inactivated, ambiguous, or neutral, and were removed from further consideration are listed in Table 5.

status	type	# attributes	list of attributes
removed	inactivated	9	blurry, narrow eyes, necklace, oval face, rosy cheeks, pointy nose, bald, mustache, pale skin
	ambiguous	2	attractive, heavy make up
	neutral	3	no beard, five-o-clock shadow, mouth slightly open
annotated	one or more disentangled single-channel controls found	16	gender, smiling, lipstick, eyeglasses, bangs, wavy hair, earrings, black hair, blond hair, sideburns, goatee, receding hairline, gray hair, suit (tie), double chin, hat
	no disentangled single-channel controls found	10	bags under eyes, big nose, high cheekbones, young, arched eyebrows, brown hair, big lips, bushy eyebrows, chubby, straight hair

Table 5. For 40 CelebA attributes, we first remove inactivated (9), ambiguous (2) or neutral (3) attributes. Our method is able to detect one or more disentangled single-channel controls for 16 out of the 26 remaining attributes.

To find attribute-specific controls we apply our method described in Section 5 on the remaining 26 attributes. We found that 16 out of the 26 remaining attributes are controllable by one or more single style channels, in a disentangled manner. Our method was not able to identify any *disentangled* single-channel controls for the other 10 attributes. All of the above attributes are listed in Table 5. The attributes for which no disentangled single-channel controls were found indeed appear to be correlated with other visual attributes (in FFHQ). For example, “bags under eyes” is correlated with eye size, “big lips” is correlated with skin color, and “high cheekbones” is correlated with smiling.

While our method could not find a disentangled single-channel control for the “young” attribute, it was able to find such controls for wrinkles, eyeglasses, and gray hair. Because all of these attributes are correlated with age, the “young” attribute can only be controlled by manipulating multiple style channels, rather than a single one.

Note that although the attribute-specific detection method of Section 5 could not detect a single-channel control for either “arched eyebrows” or “bushy eyebrows”, our locally-active detection method in Section 4 was able to find disentangled controls for these attributes: (9,30) for arched eyebrows, and (12,325) for bushy eyebrows. Thus, these could be considered as failure cases for our attribute-specific detection method.

Table 6 lists the various attributes and the single-channel controls that were detected for them. For each control we list the layer and channel number, as well as its rank by the detection method of Section 5.

region	attribute	(layer,channel,rank)	related attributes
hair	black hair	(12,479,1)	different hair color, lighting
	blond hair	(12,479,1)	gender, other hair color and style
	gray hair	(12,266,3)	glasses, gender, wrinkle and beard
	wavy hair	(6,500,1)	hair style, gender
		(8,128,2)	
		(5,92,3)	
		(6,394,7)	
		(6,323,28)	
	bangs	(3,259,1)	hair style
		(6,285,2)	
		(5,414,3)	
		(6,128,4)	
		(9,295,8)	
		(6,322,9)	
	receding hairline	(6,487,11)	hair style
		(6,504,14)	
		(5,414,1)	
		(6,322,2)	
mouth	smiling	(6,501,1)	size of face or eye
	lipstick	(15,45,1)	gender, face expression
beard	sideburns	(12,237,2)	other type of beard, gender
	goatee	(9,421,1)	other type of beard, gender
chin	double chin	(9,132,1)	size of neck, wrinkle
ear	earrings (entangled with gender)	(8,81,1)	gender, face shape
eye	glasses	(3,288,1)	gender, wrinkle and beard
		(2,175,3)	
		(3,120,4)	
		(2,97,6)	
clothes	suit (tie)	(9,441,1)	cloth style
		(8,292,2)	
		(11,358,3)	
hat	hat size	(6,223,11)	nothing change
		(5,200,7)	
overall	gender	(9,6,1)	beard,hair style

Table 6. List of attributes and the single-channel controls that were detected for them. The indices of layers and channels start from 0, while ranks start from 1.

E. Attribute Dependency

To compare the disentanglement of different image manipulation methods, we propose a general disentanglement metric for real images, which we refer to as *Attribute Dependency* (AD). Attribute Dependency measures the degree to which manipulation along a certain direction induces changes in other attributes, as measured by classifiers for those attributes. Below we share our insights regarding AD, and its implementation details. Next, we use AD to show our image manipulation method is more disentangled than two other methods (GANSpace [12], InterfaceGan [31]) in Figure 15 and Figure 16. Additionally, we further show in Figure 17 that our method changes face identity less than GANSpace and InterFaceGAN.

For a given target attribute t , we measure AD as follows. First, we sample a set of images without the target attribute t (e.g., without gray hair), and manipulate them towards the target attribute, by a certain amount measured by the change in the logit outcome Δl_t of a classifier pretrained to detect attribute t . Next, we measure the change of logit Δl_i between the original images and the manipulated ones for other attributes $\forall i \in \mathcal{A} \setminus t$, where \mathcal{A} is the set of all attributes. Each change is normalized by $\sigma(l_i)$, the standard deviation of the logit value for attribute i over a large set of generated images. We measure mean-AD, defined as $E(\frac{1}{k} \sum_{i \in \mathcal{A} \setminus t} (\frac{\Delta l_i}{\sigma(l_i)}))$, where $k = |\mathcal{A}| - 1$. Similarly, we measure max-AD, defined as $E(\max_{i \in \mathcal{A} \setminus t} (\frac{\Delta l_i}{\sigma(l_i)}))$.

E.1. Insights

To measure how much a specific attribute $i \in \mathcal{A} \setminus t$ has changed, we use a pretrained classifier for that attribute. Under normal operating mode, a binary classifier outputs a logit $l_i \in [-\infty, +\infty]$, which is then converted to a probability value in $[0, 1]$, with positive logit values yielding probabilities higher than 0.5, and negative logit values yielding probabilities lower than 0.5. However, classifiers trained on real data may be affected by entanglement present in the training data, and they are often unable to detect the presence or absence of an attribute in a disentangled manner. For example, a female face with lipstick will typically cause the classifier to output a negative logit value (indicating the presence of a lipstick), but the classifier might output a positive logit value given a face of a male with lipstick. Similarly, a gray hair classifier will output a negative logit value for a male with gray hair, but might output a positive logit value for a female with gray hair. This is demonstrated in Figures 13 and 14.

Thus, when attempting to measure the magnitude of change of an attribute, we choose not to consider the classifier’s logit sign or value; rather, we find that the change in the logit value, Δl , appears to be better correlated with an image space change in the attribute. We use the change in the logit, rather than the change in the probability because of the saturating effect of the sigmoid non-linearity that is used to convert logits to probabilities. For example, the probability produced by a lipstick classifier for a female wearing a lighter lipstick and a stronger lipstick is going to be nearly the same, while this is not the case for the logit values (see the last row of Figure 13).

Another insight is that if the manipulation strength is too high, the generated images will be unrealistic, and classifiers will give unexpected predictions. It is crucial not to use too high manipulation strengths to make sure the logits are meaningful. If the generated images are realistic, the logit is nearly a monotonic function of the strength of target attribute. And we consider the manipulation strength is a monotonic function of strength of target attribute. Therefore, we can control the amount of manipulation (measured by Δl) by searching for the corresponding manipulation strength through bisection method.

Our final insight is that the classifiers are not immune to noise. When provided with the same images with only slight texture changes in hair and background, the classifiers are supposed to output the same logits. In practice, however, the logits are slightly different. Thus, when comparing the effect of different manipulation methods on various attributes, it is necessary to ensure that the differences in the measurements are caused by the inherent differences between the methods, rather than by noise in the classifier outputs.

E.2. Implementation

We randomly generate 500K images, as our image bank, and annotate each image with 31 active attributes, same as was done in Section 3, where a negative logit corresponds to presence of the target attribute in an image. Let $\sigma(l)$ denote the standard deviation of logits over the entire image bank. For each target attribute (for example, gray hair), we rank its logit from negative to positive, and take images with 50-75% quantile as manipulation candidates, since they exhibit little to mild presence of the target attribute (without much gray hair). We don’t take images with the most positive logit (75-100% quantile) since they are less likely to result in a realistic manipulation. Candidates are manipulated toward strong attribute presence (adding gray hair, more negative logit). We set the $\Delta l_t = r\sigma(l_t)$, where $r \in \{0.25, 0.5, 0.75, 1\}$. r should not be too large to make sure that most of the manipulated images are still realistic. Then we use the bisection method to find the manipulation strength m that can generate an image with final logit $|(l_t^{final} - (l_t^{initial} - \Delta l_t))| < r_{tolerance}\sigma(l_t)$ with $r_{tolerance} = 5\%$, and ignore images that don’t converge after 20 iterations. We manually set the maximum manipulation

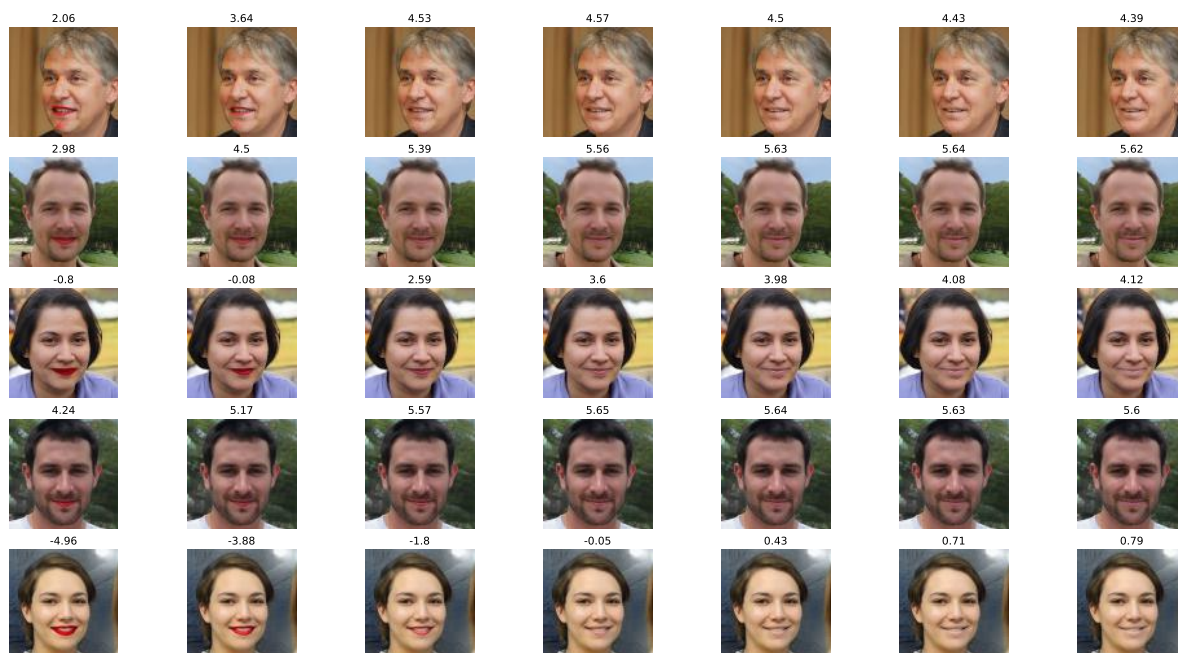


Figure 13. Logits of lipstick classifier. The strength of attribute is reduced from left to right. The classifier logit is on top of each image, increasing from left to right. Note that the logit sign is not aligned with the presence of the attribute: there are images with strong lipstick, but a positive logit.

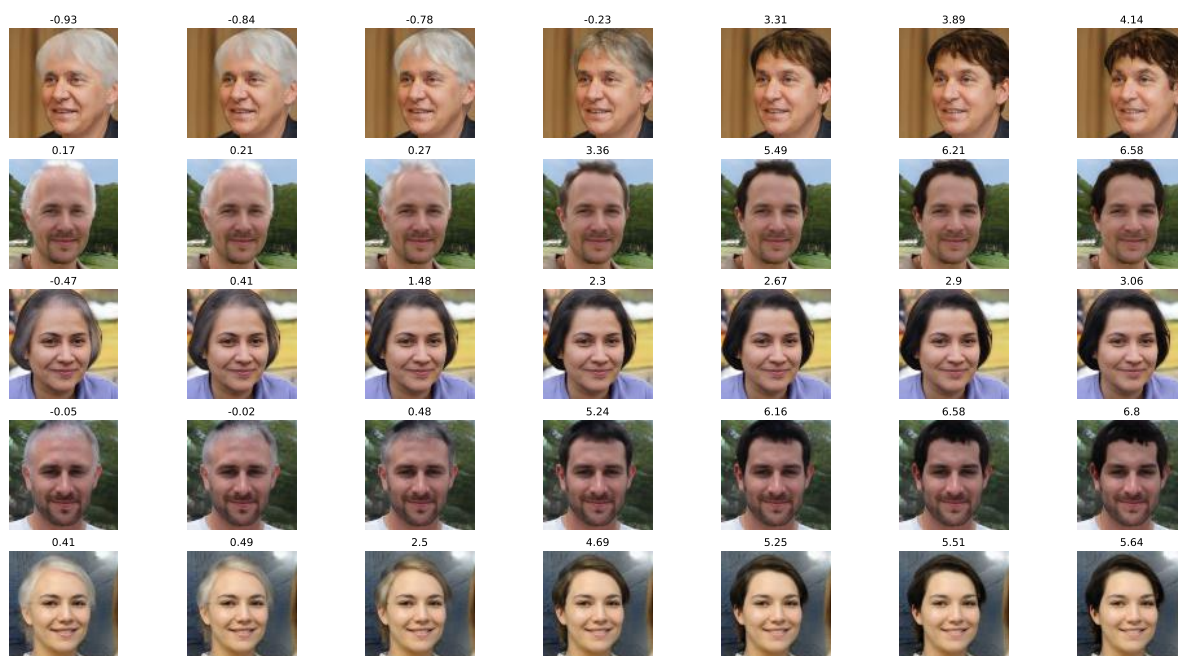


Figure 14. Logits of hair greyness classifier. The strength of attribute is reduced from left to right. The classifier logit is on top of each image, increasing from left to right. Note that the logit sign is not aligned with the presence of the attribute: there are images with strong hair greyness, but a positive logit.

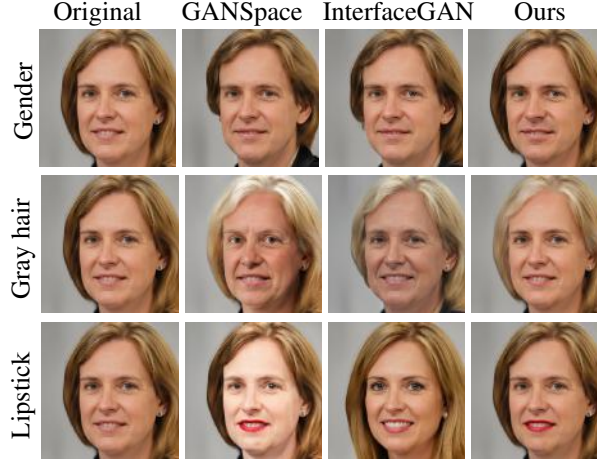


Figure 15. Comparison with state-of-the-art methods with amount of manipulation $\Delta l_t = 1.5\sigma(l_t)$. We deliberately choose a strong manipulation (1.5 instead of 1) to emphasize the differences.

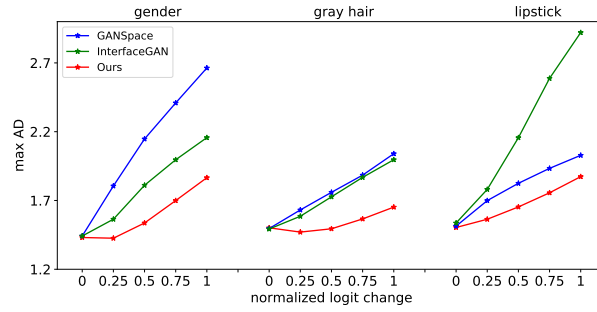


Figure 16. Max-AD vs. the degree of target attribute manipulation ($\Delta l_t / \sigma(l_t)$). Lower max-AD indicates better disentanglement.

strength m_{max} such that almost all manipulated images with manipulation strength m_{max} strongly exhibit the target attribute, but still look realistic. m_{max} is used to initiate the bisection method.

Next, we add a control group with $r = 0$ to the experiment, such that $\Delta l_t = 0\sigma(l_t)$. This group is used to represent the inherent noise of classifiers. The input images are copies of the original ones, obtained by keeping the latent code s unchanged, but changing the noise inputs at different layers. They are essentially identical to the original images, with subtle differences in hair, skin, and background.

Finally, we use the same 3K images with identifiable m for all candidate manipulation methods to calculate mean-AD and max-AD. The mean-AD for the three methods (GANSpace, InterFaceGAN, and ours) for three attributes (gender, gray hair, and lipstick) are plotted in Figure 8, and the max-AD in Figure 16. Figures 6 and 15 show a qualitative comparison between the manipulations produced by the three methods. Note that, like in Figure 6, the *Lipstick* manipulation by InterFaceGAN significantly changes the identity of the person, and the *Gray hair* manipulation adds wrinkles. GANSpace manipulations also exhibit some entanglement (*Lipstick* affects face lightness, *Gray hair* ages the rest of the face). In contrast, our approach appears to affect only the target attribute. Our *Gender* manipulation, for example, does not affect the hair style, and minimally changes the face, yet the gender unmistakably changes.

E.3. Identity change

In addition to AD, we use another metric (identity change) to compare our method against GANSpace and InterFaceGAN. Specifically, we use FaceNet [29], which is a standard network for measuring identity change. We use the official implementation, which first detects faces, then crops faces out, obtains an embedding from the last layer, and calculates the Euclidean norm between the embedding of the original images and that of the manipulated ones. As shown in Figure 17, manipulations done with our method change the identity less than manipulations of GANSpace or InterFaceGAN.

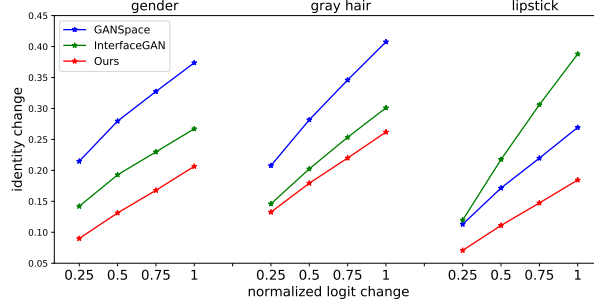


Figure 17. Identity change vs. the degree of target attribute manipulation ($\Delta l_t / \sigma(l_t)$). Lower identity changes indicates that the manipulation is better disentanglement from the identity.

F. Manipulation of real images

To manipulate real images, it is necessary to first invert them into latent codes. Through latent optimization [18], we observe that the reconstruction quality is the highest when optimizing in \mathcal{S} , followed by $\mathcal{W}+$, and is the lowest for \mathcal{W} , as demonstrated in Figure 18. However, the naturalness of subsequent manipulation is the best when the latent optimization is done in \mathcal{W} , followed by $\mathcal{W}+$, and the worst for \mathcal{S} , as shown in Figure 19. Through training a latent embedding encoder and using the embedding produced by the encoder as the initial point for a few iterations of latent optimization, we obtain both good reconstruction and natural manipulation. We demonstrate manipulation of real images in Figure 20 (for real images from the FFHQ dataset) and in Figure 21 for images from the CelebA-HQ [20] dataset.

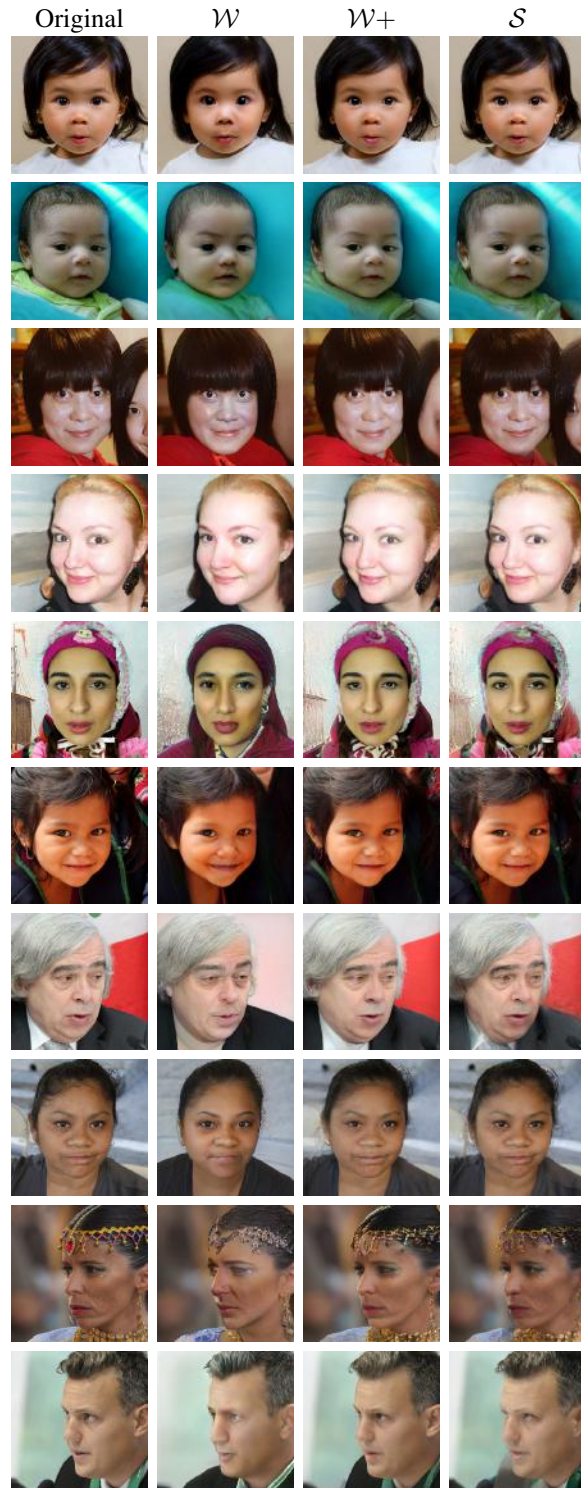


Figure 18. Inversion via latent optimization in \mathcal{W} , $\mathcal{W}+$, \mathcal{S} . It may be easily seen that the reconstruction is least accurate for \mathcal{W} , more accurate for $\mathcal{W}+$, and is best for \mathcal{S}

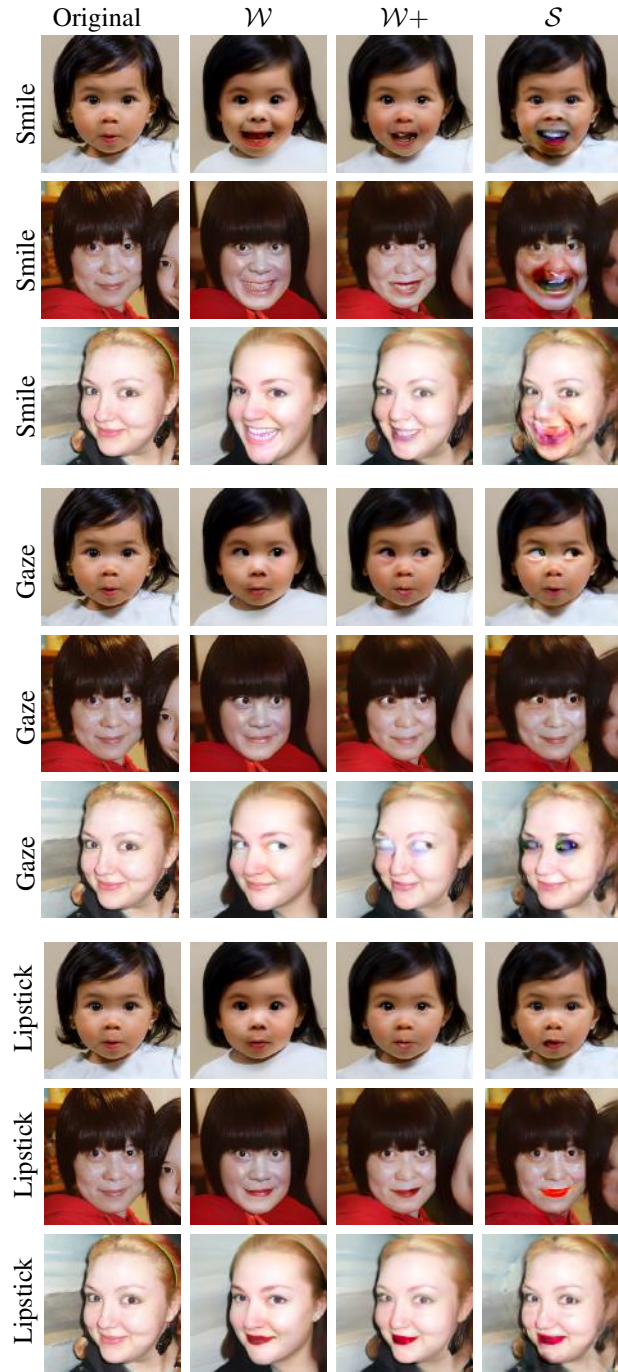


Figure 19. Manipulation of style codes obtained by latent optimization in \mathcal{W} , $\mathcal{W}+$, and \mathcal{S} spaces. The exact same manipulation is applied in each row. It may be seen that manipulation of codes optimized in \mathcal{S} produces significant artifacts, while manipulation on codes optimized in \mathcal{W} produce more realistic results.

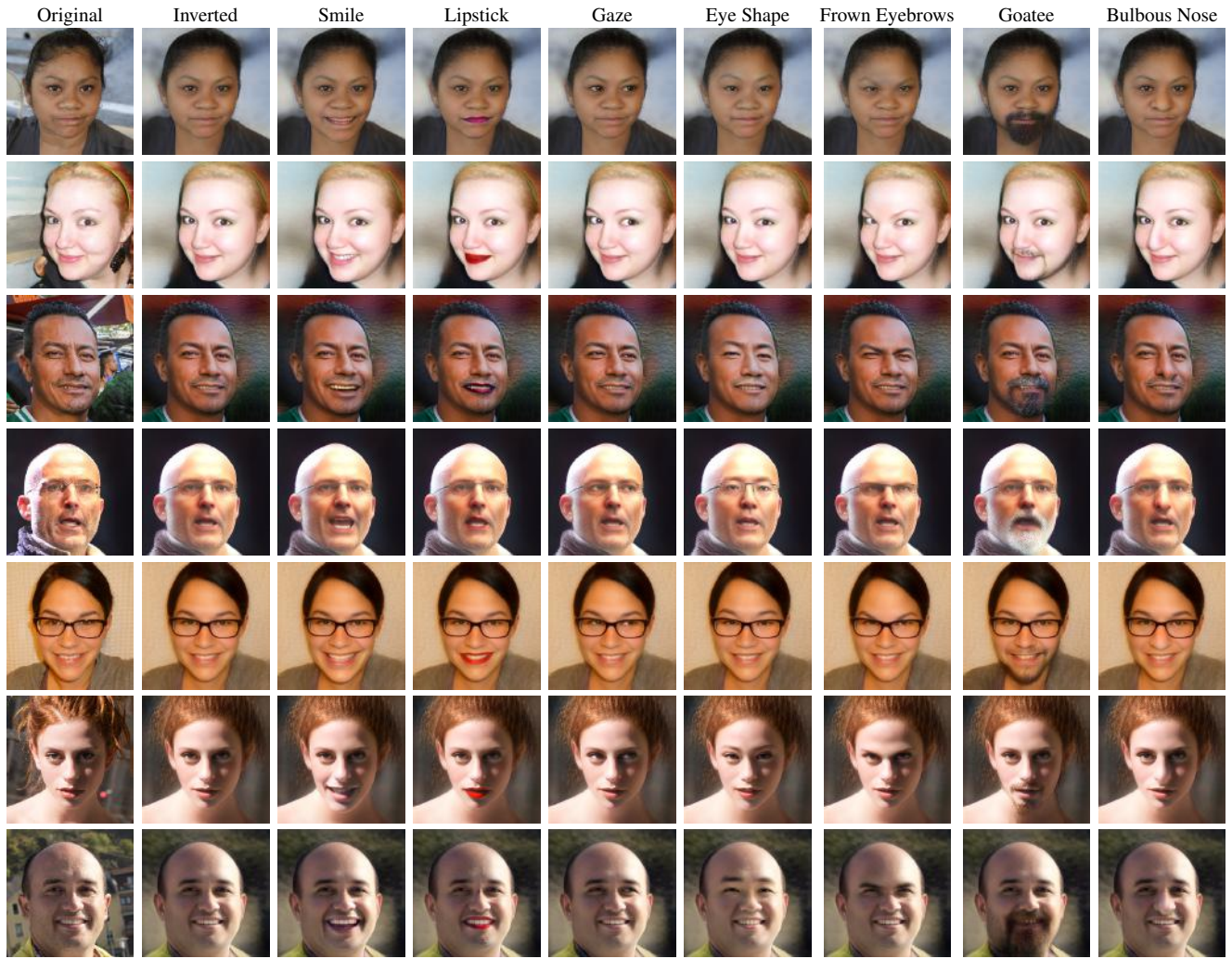


Figure 20. Manipulation of real images using encoder-based inversion. Original images are from FFHQ, and were not part of the encoder’s training set.

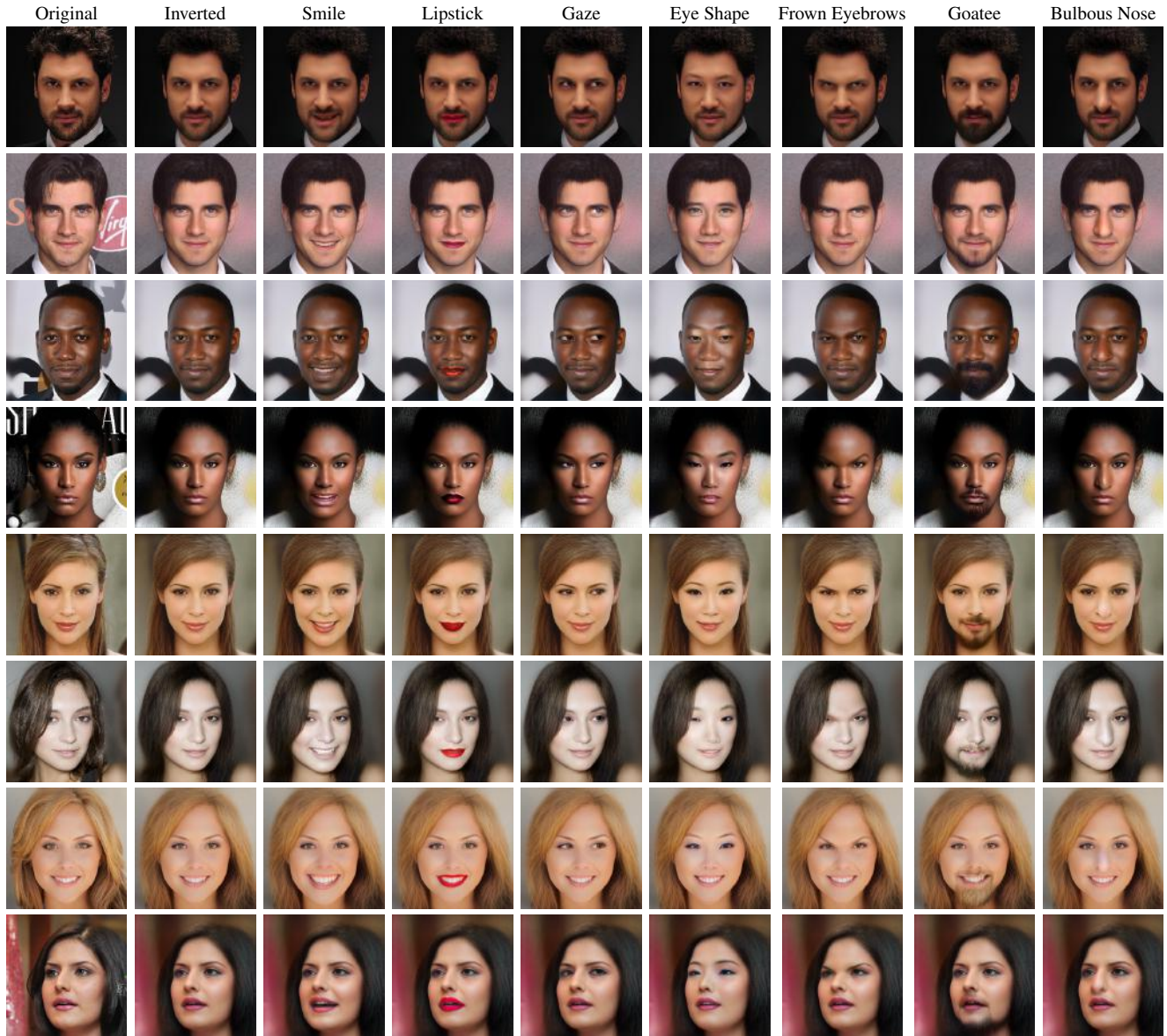


Figure 21. Manipulation of real images using encoder-based inversion. Original images are from CelebA-HQ, which were not part of the encoder’s training set, and not part of the GAN training set (the StyleGAN2 model was trained on the FFHQ dataset).